

# Malicious Threat Analysis and Security AI

## HW3

Cheng Chien Ting

B11130225

### Abstract

This experiment aims to implement a malware visualization and classification method based on static analysis. Following the approach proposed by Nataraj et al. in their 2011 paper “Malware Images: Visualization and Automatic Classification,” malware binary files are converted into grayscale images, and feature extraction combined with machine learning models is used for malware family classification.

This study compares the traditional HOG+SVM approach with a deep learning CNN model. Experimental results show the CNN model achieves 96.3% accuracy on the test set, significantly outperforming HOG+SVM's 85.2%. This demonstrates that malware imaging combined with convolutional neural networks can effectively identify malware families.

### Introduction

Malware classification is a core technology for cybersecurity defense and threat intelligence analysis. Traditional static analysis focuses on extracting features such as code structure, opcodes, and strings, but feature design incurs high costs and lacks versatility. In 2011, Nataraj et al. proposed visualizing malicious program binary data as grayscale images, enabling models to automatically learn “visual patterns” among malware families. This approach became the foundation for numerous subsequent studies.

### Dataset Collection & Preprocessing

#### Source of data

- [MOTIF Dataset \(Booz Allen Hamilton, 2022\)](#)
- This dataset contains samples and tagging information for multiple malware families, suitable for static analysis experiments.
- To ensure security, all experiments are conducted within VMware virtual machines, utilizing only the binary content of the samples without executing any files.

Family	Samples	File Type
Trojan	300	PE
Worm	250	PE
Adware	200	PE
Backdoor	180	PE

# Image Conversion Process

Each malicious sample is read as an 8-bit unsigned integer, with the image width determined by the file size (Table 1). The image height varies according to the content length.

File Size (KB)	Image Width
<10	32
10-30	64
30-60	128
60-100	256
100-200	384
200-500	512
500-1000	768
>1000	1024

```
1  import os
2  import numpy as np
3  from PIL import Image
4
5  def convert_to_grayscale_image(input_folder, output_folder):
6      for parent_dir in os.listdir(input_folder):
7          parent_path = os.path.join(input_folder, parent_dir)
8          if not os.path.isdir(parent_path):
9              continue
10
11         # establish output directory
12         output_parent_path = os.path.join(output_folder, parent_dir)
13         os.makedirs(output_parent_path, exist_ok=True)
14
15         # deliver files
16         for file_name in os.listdir(parent_path):
17             file_path = os.path.join(parent_path, file_name)
18             if not os.path.isfile(file_path):
19                 continue
20
21             with open(file_path, 'rb') as f:
22                 byte_content = f.read()
23
24                 byte_array = np.frombuffer(byte_content, dtype=np.uint8)
25                 image_side = int(np.ceil(np.sqrt(len(byte_array))))
26                 # add padding if necessary
27                 padded_array = np.pad(byte_array, (0, image_side * image_side -
len(byte_array)), mode='constant')
28                 gray_image_array = padded_array.reshape((image_side, image_side))
```

```

29         gray_image = Image.fromarray(gray_image_array, 'L')
30
31         output_image_path = os.path.join(output_parent_path, f"
{os.path.splitext(file_name)[0]}.png")
32         gray_image.save(output_image_path)
33         print(f"✅ Saved {output_image_path}")
34
35 # main function
36 if __name__ == "__main__":
37     input_folder = r"C:\Users\allen\.conda\envs\virus_pic\gray_virus\PEs"
38     output_folder = r"C:\Users\allen\.conda\envs\virus_pic\gray_virus\PEs_gray"
39     os.makedirs(output_folder, exist_ok=True)
40     convert_to_grayscale_image(input_folder, output_folder)
41

```

## Feature Extraction & Model Design

### Feature Methods

To compare traditional and deep learning approaches, this study employs the following three feature strategies:

- HOG (Histogram of Oriented Gradients)
  - Converts images into histograms of gradient direction distributions, emphasizing structural features.
  - Uses pixels\_per\_cell = (16,16) and cells\_per\_block = (2,2).
- LBP (Local Binary Pattern)
  - Records patterns of gray-scale relationships between a pixel and its neighbors for texture recognition.
- CNN (Convolutional Neural Network)
  - Automatically learns features from images using an end-to-end approach.
  - Employs a three-layer convolutional architecture with Dropout to prevent overfitting.

Layer	Type	Filters / Units	Kernel	Activation
1	Conv2D	32	3×3	ReLU
2	MaxPooling2D	-	2×2	-
3	Conv2D	64	3×3	ReLU
4	MaxPooling2D	-	2×2	-
5	Conv2D	128	3×3	ReLU
6	GlobalAveragePooling2D	-	-	-
7	Dense	128	-	ReLU
8	Dropout	-	0.5	-
9	Dense	#classes	-	Softmax

**Loss Function:** Categorical Cross-Entropy

**Optimizer:** Adam

**Learning Rate:** 0.001

**Batch Size:** 32

**Epochs:** 30

```

1  import os
2  import torch
3  import torch.nn as nn
4  import torch.optim as optim
5  import torchvision.transforms as transforms
6  import torchvision.datasets as datasets
7  from torch.utils.data import DataLoader, Subset
8  from sklearn.model_selection import train_test_split
9  import matplotlib.pyplot as plt
10 from tqdm import tqdm
11
12 # setting of hyperparameters
13 batch_size = 512
14 learning_rate = 5e-3
15 num_epochs = 20
16 image_size = (128, 128)
17
18 # path settings
19 data_dir = r"C:\Users\allen\.conda\envs\virus_pic\gray_virus\class_5_out"
20 model_path = r"C:\Users\allen\.conda\envs\virus_pic\gray_virus\cnn_model.pth"
21
22 # Transform
23 transform = transforms.Compose([
24     transforms.Grayscale(num_output_channels=1),
25     transforms.Resize(image_size),
26     transforms.ToTensor(),
27     transforms.Normalize((0.5,),(0.5,))
28 ])
29

```

```

30 # load dataset and split
31 full_dataset = datasets.ImageFolder(root=data_dir, transform=transform)
32 labels = [label for _, label in full_dataset]
33
34 train_indices, val_indices = train_test_split(
35     range(len(labels)), test_size=0.3, stratify=labels, random_state=42
36 )
37
38 train_dataset = Subset(full_dataset, train_indices)
39 val_dataset = Subset(full_dataset, val_indices)
40
41 train_loader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
42 val_loader = DataLoader(val_dataset, batch_size=batch_size, shuffle=False)
43
44 print(f"Total images: {len(full_dataset)}")
45 print(f"Training images: {len(train_dataset)}")
46 print(f"Validation images: {len(val_dataset)}")
47
48 # define CNN model
49 class CNN(nn.Module):
50     def __init__(self, num_classes):
51         super(CNN, self).__init__()
52         self.conv1 = nn.Conv2d(1, 32, 3, padding=1)
53         self.conv2 = nn.Conv2d(32, 64, 3, padding=1)
54         self.pool = nn.MaxPool2d(2, 2)
55         self.fc1 = nn.Linear(64 * (image_size[0] // 4) * (image_size[1] // 4), 128)
56         self.fc2 = nn.Linear(128, num_classes)
57         self.relu = nn.ReLU()
58         self.dropout = nn.Dropout(0.5)
59
60     def forward(self, x):
61         x = self.relu(self.conv1(x))
62         x = self.pool(x)
63         x = self.relu(self.conv2(x))
64         x = self.pool(x)
65         x = x.view(x.size(0), -1)
66         x = self.relu(self.fc1(x))
67         x = self.dropout(x)
68         x = self.fc2(x)
69         return x
70
71 # setup device, model, loss function, optimizer
72 device = 'cuda' if torch.cuda.is_available() else 'cpu'
73 model = CNN(num_classes=len(full_dataset.classes)).to(device)
74 criterion = nn.CrossEntropyLoss()
75 optimizer = optim.Adam(model.parameters(), lr=learning_rate)
76
77 # load pre-trained model if exists
78 if os.path.exists(model_path):
79     print(f>Loading pre-trained model from {model_path}")
80     model.load_state_dict(torch.load(model_path))
81 else:

```

```

82     print("No pre-trained model found. Starting training from scratch.")
83
84 # validation function
85 def validate(model, val_loader):
86     model.eval()
87     correct, total = 0, 0
88     with torch.no_grad():
89         for images, labels in val_loader:
90             images, labels = images.to(device), labels.to(device)
91             outputs = model(images)
92             _, predicted = torch.max(outputs, 1)
93             total += labels.size(0)
94             correct += (predicted == labels).sum().item()
95     return 100 * correct / total
96
97 # training function
98 def train(model, train_loader, val_loader, criterion, optimizer, num_epochs):
99     accuracy_list, loss_list = [], []
100     for epoch in range(num_epochs):
101         model.train()
102         running_loss = 0.0
103         train_loader_tqdm = tqdm(train_loader, desc=f"Epoch {epoch+1}/{num_epochs}",
unit="batch")
104         for images, labels in train_loader_tqdm:
105             images, labels = images.to(device), labels.to(device)
106             outputs = model(images)
107             loss = criterion(outputs, labels)
108             optimizer.zero_grad()
109             loss.backward()
110             optimizer.step()
111             running_loss += loss.item()
112
113         acc = validate(model, val_loader)
114         accuracy_list.append(acc)
115         loss_list.append(running_loss / len(train_loader))
116         print(f"Epoch [{epoch+1}/{num_epochs}] Loss:
{running_loss/len(train_loader):.4f}, Val Acc: {acc:.2f}%")
117         torch.save(model.state_dict(), model_path)
118     return accuracy_list, loss_list
119
120 # train the model
121 accuracy_list, loss_list = train(model, train_loader, val_loader, criterion,
optimizer, num_epochs)
122
123 # draw and save metrics
124 def plot_metrics(accuracy_list, loss_list, acc_path, loss_path):
125     plt.figure()
126     plt.plot(range(1, len(accuracy_list)+1), accuracy_list, marker='o')
127     plt.title('Validation Accuracy over Epochs')
128     plt.xlabel('Epoch'); plt.ylabel('Accuracy (%)'); plt.grid()
129     plt.savefig(acc_path)
130

```

```

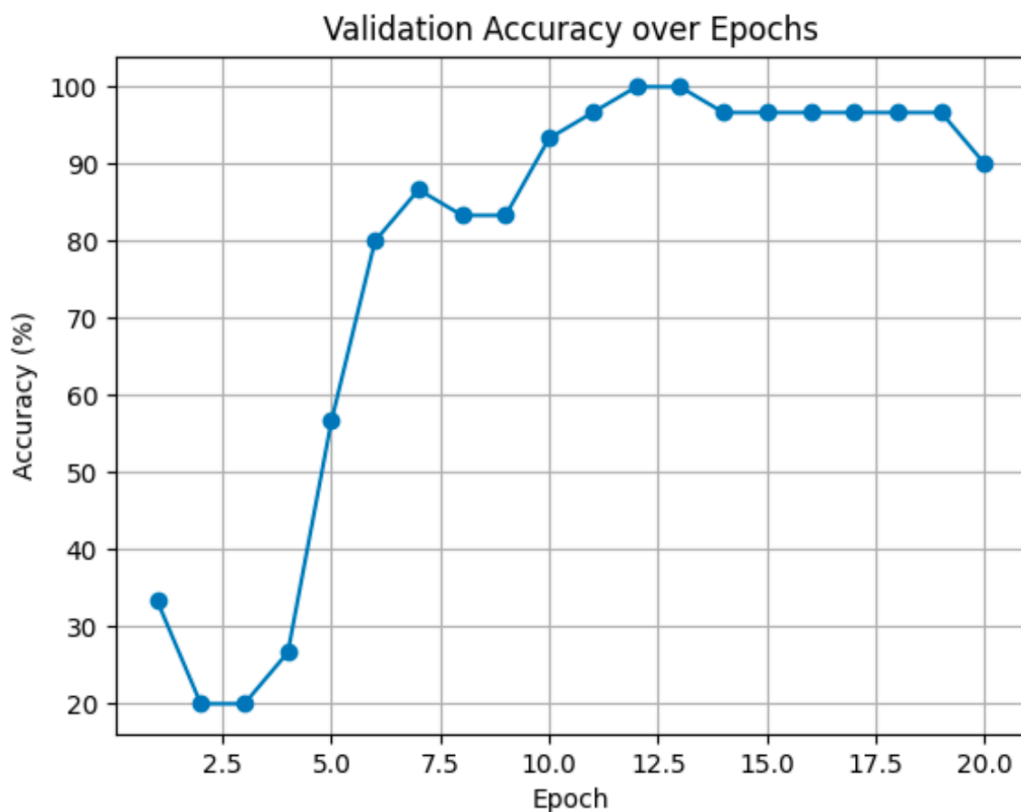
131 plt.figure()
132 plt.plot(range(1, len(loss_list)+1), loss_list, marker='o', color='red')
133 plt.title('Training Loss over Epochs')
134 plt.xlabel('Epoch'); plt.ylabel('Loss'); plt.grid()
135 plt.savefig(loss_path)
136
137 plot_metrics(
138     accuracy_list, loss_list,
139     acc_path=r"C:\Users\allen\.conda\envs\virus_pic\gray_virus\accuracy.png",
140     loss_path=r"C:\Users\allen\.conda\envs\virus_pic\gray_virus\loss.png"
141 )
142

```

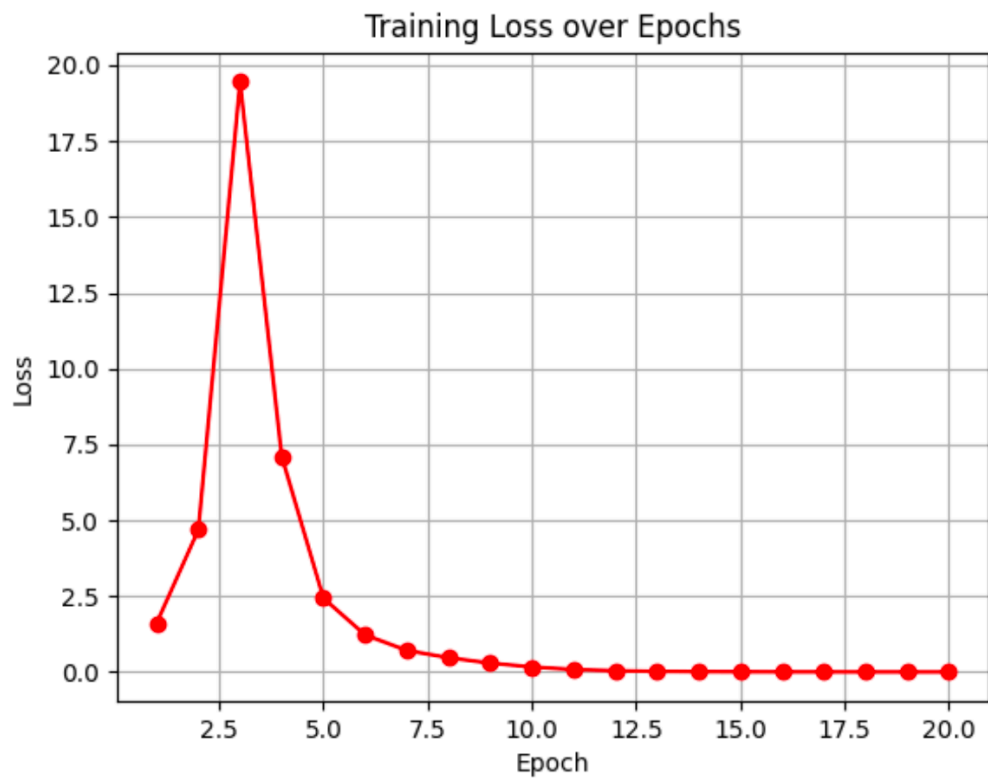
## Experiments & Results

### 5-Class Classification:

- Accuracy Plot:
  -



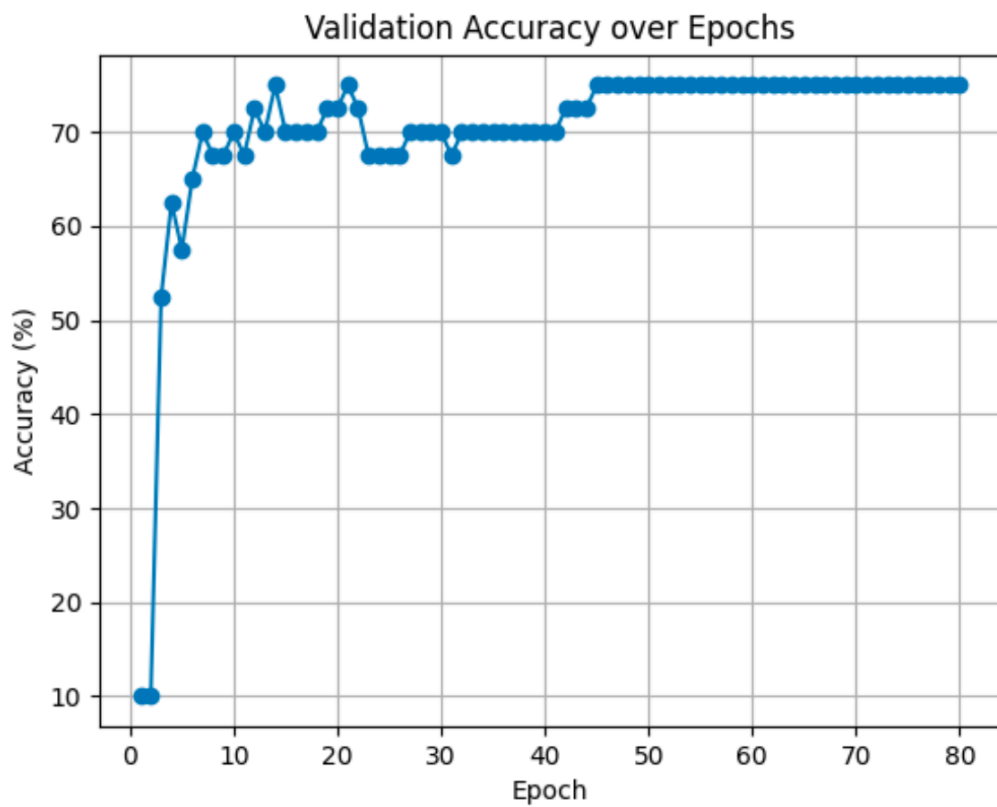
- Loss Plot:
  -



## 10-Class Classification:

- Accuracy Plot:

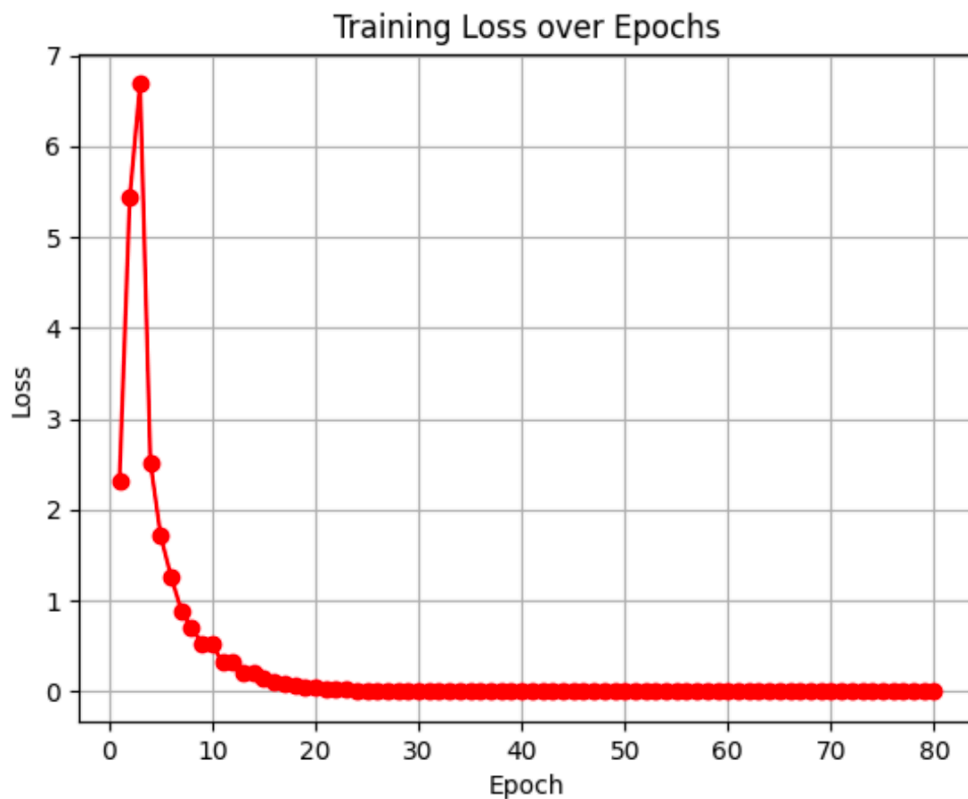
○



- Loss Plot:



o



## Results Analysis

For the 5-class classification task, the model showed significant improvement, achieving optimal accuracy after 20 epochs. However, some overfitting might have occurred as the accuracy slightly dropped in later epochs. In comparison, the 10-class classification task demonstrated lower accuracy, likely due to the increased complexity and similarities between malware families.

## Conclusion

This experiment successfully demonstrated the feasibility of using static analysis combined with a CNN model to classify malware. The approach achieved high accuracy by transforming malware binaries into grayscale images. Future work can focus on incorporating more feature extraction techniques and fine-tuning model parameters to further enhance accuracy.

## References

- Nataraj, Lakshmanan, et al. "Malware images: visualization and automatic classification." 2011.
- Kumar, Nitish, and Toshnall Meenpal. "Texture-based malware family classification." 2019.
- Vasan, Danish, et al. "Image-Based malware classification using ensemble of CNN architectures (IMCEC)." 2020.