



Sri Lanka Institute of Information Technology

# Deep Learning for Malware Detection: A CNN-Based Approach

## Individual Assignment

IE4032 - Information Warfare

Submitted by:

Student Registration Number	Student Name
IT20128272	Rathnayake .R .M .K .G

3<sup>rd</sup> of November, 2023

Date of submission

## Acknowledgment

I would like to express my sincere gratitude to my esteemed lecturer, Dr. Lakmal Rupasinghe, for their invaluable guidance and mentorship throughout the course of this project. Their expertise and unwavering support have been instrumental in my journey towards understanding and implementing deep learning for malware analysis.

I also extend my heartfelt acknowledgment to the faculty members of the Sri Lanka Institute of Information Technology's Computing department. Their assistance and the nurturing learning environment they provide have played a pivotal role in advancing my studies. Their insights and guidance have served as a guiding light, enabling me to navigate through the complexities of this research.

Lastly, I offer my deepest thanks to my families and friends. Your unflagging encouragement and unwavering belief in my abilities have been the driving force behind this project endeavor. Your unwavering support has been a wellspring of inspiration, motivating me to strive for excellence in every aspect of my work. I am truly grateful for your unwavering support and encouragement throughout this journey.

## Declaration

I declare that this is my own work and this project does not incorporate without acknowledgement any material previously submitted for a degree or diploma in any other university or Institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text. Also, I hereby grant to Sri Lanka Institute of Information Technology, the nonexclusive right to reproduce and distribute my dissertation, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works (such as articles or books).

Name	Student ID	Signature
Rathnayake. R .M .K .G	IT20128272	

## 1. Executive Summary

In response to the ever-evolving cybersecurity landscape, this report details the development of a cutting-edge malware detection system powered by Convolutional Neural Networks (CNNs) designed to analyze malware images. The primary objective of this project was to create a deep learning model capable of effectively distinguishing between benign and malware images, with a strong emphasis on practical real-world application.

Our foremost achievement lies in the creation of an exceptionally efficient malware detection model. The CNN model exhibited remarkable accuracy, proving its proficiency in classifying images as either malware or benign. This accomplishment was a result of diligent model training and evaluation, highlighting its potential as a robust solution for malware detection.

Another key aspect of our success was in the versatility and thoroughness of our data preprocessing procedures. These encompassed a comprehensive series of steps, including data conversion, normalization, and feature selection. These measures significantly optimized the model's performance by ensuring the use of high-quality input data, a cornerstone in the process of crafting a dependable malware detection system.

Furthermore, we introduced a practical use case scenario where our CNN model serves as the foundation for a software product. This software enables users to inspect .exe files, a common carrier of malware, and promptly receive feedback on their malware status. The system architecture was thoughtfully designed to enhance cybersecurity practices, thereby contributing to the safety of computing devices.

In terms of future directions, we underscore the importance of ongoing improvements. This includes continuous model enhancement through training and fine-tuning to adapt to emerging malware threats, as well as expanding support to a broader range of file types to increase the software's versatility. Enhancements to the user interface, additional features, and user-friendly improvements are pivotal in providing a seamless user experience.

We are eager to explore opportunities for large-scale deployment, serving the needs of both personal users and enterprises seeking robust malware detection solutions. Integration with existing cybersecurity systems offers the potential to create a comprehensive defense strategy, while a transition to real-time analysis ensures immediate protection against malware threats. In conclusion, this project represents a significant stride in leveraging deep learning to enhance cybersecurity practices. Our CNN model, in tandem with the envisioned software product, contributes to a safer and more secure digital environment, with ample opportunities for further advancements in the ever-evolving field of cybersecurity.

## Table of Contents

Acknowledgment.....	2
Declaration .....	3
1. Executive Summary .....	4
2. Technical Review .....	7
2.1. Introduction.....	7
2.2. Methodology .....	8
2.2.1. CNN Model .....	8
2.2.2. The Dataset.....	8
2.2.3. Model Training .....	9
2.2.4. Google Colab .....	10
2.2.5. Libraries .....	10
2.2.6. Data Loading and Exploration .....	13
2.2.7. Data Preprocessing.....	15
2.2.8. Data Analysis and Feature Selection .....	16
2.2.9. Data Splitting and Normalization .....	21
2.2.10. Data Shape and Dimension Adjustment .....	22
2.2.11. CNN Model Definition .....	24
2.2.12. Model Compilation and Training.....	26
2.2.13. Model Evaluation.....	28
2.3. Use Case Scenario: Malware Detection Software.....	30
2.3.1. Scenario Description:.....	30
2.3.2. User Interaction:.....	30
2.3.3. System Architecture: .....	31
3. Conclusion and Future Work.....	32
3.1. Conclusion .....	32
3.2. Future Work .....	33
References.....	34

## 2. Technical Review

### 2.1. Introduction

In the ever-evolving landscape of cybersecurity, the relentless ingenuity of malicious actors poses an ever-growing threat to digital systems and personal privacy. Malware, a portmanteau of "malicious software," represents a significant facet of these cyber threats. Malware encompasses a diverse array of programs or code explicitly designed to compromise the integrity, confidentiality, or availability of computer systems. From viruses and worms to Trojans and ransomware, malware manifests in countless forms, continually adapting to evade detection and wreak havoc on unsuspecting victims.

Traditional malware detection mechanisms, rooted in signature-based methods, have long served as the first line of defense against such threats. These methods rely on known patterns or signatures of malicious code to identify and quarantine malware. However, as cybercriminals become increasingly adept at crafting new, polymorphic, and obfuscated variants of malware, the limitations of these conventional approaches become strikingly apparent. The need for innovative and adaptive solutions is paramount.

Enter deep learning, a subset of machine learning, known for its ability to learn intricate patterns and representations directly from data. Deep learning models, particularly Convolutional Neural Networks (CNNs), have gained significant prominence across various domains, from image recognition to natural language processing. This technology's prowess in feature extraction and classification has opened up a promising frontier for tackling the complex problem of malware detection and analysis.

In this report, we delve into the application of deep learning, specifically CNNs, for the purpose of malware image detection and analysis. By harnessing the power of neural networks and the hierarchical features they can uncover within images, we aim to significantly enhance the accuracy and robustness of malware detection. Through this exploration, we seek to answer critical questions, including how deep learning can revolutionize traditional malware detection mechanisms, and what unique advantages it offers in the ongoing battle against cyber threats.

The overarching purpose of this project is to develop and evaluate a deep learning model using CNNs for the detection of malware within images. By doing so, we aim to contribute to the broader field of cybersecurity by enhancing the efficiency and adaptability of malware detection, ultimately reinforcing the resilience of digital systems and data privacy in an era where cyber threats are increasingly sophisticated and insidious.

In the upcoming sections, we will explore the methodology, practical applications, findings, and in-depth discussions that collectively form our pursuit of strengthening malware detection using deep learning approaches.

## 2.2. Methodology

### 2.2.1. CNN Model

To undertake the task of malware image analysis, we employed a Convolutional Neural Network (CNN) as the cornerstone of our deep learning model. CNNs are particularly well-suited for image-related tasks due to their inherent ability to automatically learn and extract meaningful features from images. These networks consist of convolutional layers, pooling layers, and fully connected layers that work in tandem to progressively abstract high-level information from the input images.

Our CNN architecture was implemented using Keras, a high-level neural networks API that operates on top of TensorFlow. We configured our model with several convolutional layers, followed by max-pooling layers to reduce spatial dimensions, and fully connected layers for classification. The choice of CNN architecture was guided by a balance between complexity and computational feasibility, with the aim of capturing important features while maintaining model efficiency.

### 2.2.2. The Dataset

To facilitate the training and evaluation of our CNN model, we employed a carefully curated dataset. This dataset includes various attributes and features relevant to malware images, such as hash values, millisecond timestamps, classification labels, system state information, and several performance metrics. Notably, this dataset was sourced from



Kaggle, a reputable platform for machine learning datasets, ensuring the quality and diversity of the data used in our project.

The dataset attributes are as follows:

- **Hash:** Unique identifier for each malware sample.
- **Millisecond:** Timestamp providing temporal context.
- **Classification:** Labels indicating the malware type or category.
- **State:** Information on the system's state during the malware detection.
- **Usage Counter:** Counter for tracking the usage of the system.
- **Prioritization Metrics:** Includes attributes like prio, static\_prio, normal\_prio, and policy.
- **Virtual Memory Management:** Attributes like vm\_pgoff, vm\_truncate\_count, task\_size, cached\_hole\_size, free\_area\_cache.
- **Memory Metrics:** mm\_users, map\_count, hiwater\_rss, total\_vm, shared\_vm, exec\_vm, reserved\_vm, nr\_ptes.
- **Process Metrics:** end\_data, last\_interval, nvcsw, nivcsw, minflt, majflt, fs\_excl\_counter.
- **Lock:** Lock information related to resource access.
- **CPU Time:** User time (utime), system time (stime), guest time (gtime), and cumulative guest time (cgtime).
- **Signal Information:** Information regarding signal handling and context switching (signal\_nvcsw).

### 2.2.3. Model Training

The core of our approach lies in the training of the CNN model. We employed a Keras Sequential model for this purpose. Training involved feeding the preprocessed dataset into the model and optimizing the model's parameters through backpropagation. We used standard training practices, including batch training, data augmentation, and dropout layers, to improve the model's generalization capability and minimize overfitting.

The dataset was split into training and validation subsets to monitor the model's performance during training. Metrics such as loss and accuracy were tracked to gauge the

model's progress. The model was trained over multiple epochs to refine its feature extraction and classification abilities.

The effectiveness of our CNN model was evaluated using various performance metrics and these metrics provide insights into the model's ability to correctly classify malware images and distinguish them from benign images.

By following this methodology, we have laid the foundation for the successful development and training of our deep learning model for malware image detection. The subsequent sections of this report will delve into the results and implications of our approach.

#### 2.2.4. Google Colab

The implementation of our deep learning model for malware image analysis was carried out on Google Colab, a cloud-based Jupyter notebook environment provided by Google. Google Colab offers several key advantages, including access to powerful GPUs and TPUs, which significantly expedite the training of complex deep learning models, such as the CNN used in our project. This cloud-based platform allowed us to leverage Google's computing resources, ensuring that our model could be trained efficiently and without the constraints of local hardware limitations. Furthermore, Google Colab's seamless integration with Google Drive simplified data storage and retrieval, making it a convenient and practical choice for our project's implementation. The collaborative nature of Google Colab also facilitated teamwork and sharing of code and results among project members. This cloud-based approach to model development underscores the adaptability and accessibility of deep learning in addressing real-world challenges like malware detection and analysis.

#### 2.2.5. Libraries

The successful implementation of our deep learning model for malware image analysis relied on a robust set of libraries and frameworks that empowered us to work with data, build, train, and evaluate our Convolutional Neural Network (CNN) model. The following libraries played pivotal roles in our project:

```
[ ] import os
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import tensorflow as tf
import keras

from keras.layers import Conv2D, MaxPooling2D, Flatten, Dense
from keras.optimizers import Adam
from keras.models import Sequential
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
```

### **Python:**

The core programming language used for the project, providing a versatile and expressive foundation for data manipulation, model development, and evaluation.

### **Pandas:**

This data manipulation library enabled us to handle and preprocess the dataset efficiently. We utilized Pandas for tasks such as data loading, exploration, and transformation.

### **NumPy:**

NumPy, a fundamental library for numerical operations, facilitated data manipulation and array handling, which is essential for image processing and model training.

### **Matplotlib and Seaborn:**

These libraries were essential for data visualization. We used them to create insightful graphs and plots, helping us understand the dataset and model performance.

### **TensorFlow and Keras:**

TensorFlow served as the deep learning framework underpinning our project, while Keras, a high-level API, streamlined the process of building, training, and evaluating our CNN

model. Keras provided a user-friendly interface for defining neural network architectures, configuring training parameters, and tracking model performance.

### **Scikit-Learn:**

Although primarily known for machine learning, Scikit-Learn's utility extended to data preprocessing and splitting tasks. We used it for splitting our dataset into training and validation subsets and for standardizing image pixel values.

The CNN model architecture was constructed using Keras, which leverages TensorFlow as its backend. This allowed us to define the model's layers, loss functions, and optimization algorithms with ease. Our CNN model consisted of Conv2D layers for feature extraction, MaxPooling2D layers for spatial down-sampling, and fully connected Dense layers for classification.

Moreover, we employed the Adam optimizer to facilitate model training by adjusting weights and biases during the backpropagation process. This optimization algorithm played a critical role in enhancing the convergence and efficiency of our model.

In addition, the `train_test_split` function from Scikit-Learn enabled the division of our dataset into training and validation sets, facilitating model evaluation and performance tracking. Furthermore, we used the `StandardScaler` for image pixel value standardization to ensure uniformity in data distribution.

The judicious selection and utilization of these libraries and tools were instrumental in achieving the project's goals, from preprocessing the dataset to constructing, training, and evaluating our deep learning model for malware image analysis.

### 2.2.6. Data Loading and Exploration

In our pursuit of building an effective malware image analysis model, the initial stages involved data acquisition and exploration. The following code snippets showcase the essential steps taken in this process:

```
[ ] for dirname, _, filenames in os.walk('/content/drive/MyDrive/IW_Assignment'):
    for filename in filenames:
        print(os.path.join(dirname, filename))

/content/drive/MyDrive/IW_Assignment/Malware dataset.csv
```

This code demonstrates the use of the **os.walk** function to traverse a specified directory. In our case, we navigated to the directory where our dataset resided. This code allowed us to inspect and list the files present in the designated location, confirming that we were operating on the correct dataset.

```
[ ] # Load the malware dataset
raw_data = pd.read_csv("/content/drive/MyDrive/IW_Assignment/Malware dataset.csv")
```

```
[ ] # Print the column names
print(raw_data.columns)

Index(['hash', 'millisecond', 'classification', 'state', 'usage_counter',
      'prio', 'static_prio', 'normal_prio', 'policy', 'vm_pgoff',
      'vm_truncate_count', 'task_size', 'cached_hole_size', 'free_area_cache',
      'mm_users', 'map_count', 'hiwater_rss', 'total_vm', 'shared_vm',
      'exec_vm', 'reserved_vm', 'nr_ptes', 'end_data', 'last_interval',
      'nvcsw', 'nivcsw', 'minflt', 'majflt', 'fs_excl_counter', 'lock',
      'utime', 'stime', 'gtime', 'cftime', 'signal_nvcsw'],
      dtype='object')
```

Here, we employed the Pandas library to load the dataset from the specified CSV file. Once the data was loaded, we printed the column names to gain insight into the dataset's structure and variables. This step was vital in understanding the attributes we were working with, such as hash values, timestamps, classification labels, and various other metadata.

Additionally, we used the `.info()` method to examine the data types of each column and verify that the data was loaded correctly. This exploration step was crucial in ascertaining the integrity and suitability of the dataset for further processing and model development.

These preliminary actions laid the foundation for our subsequent data preprocessing and model building stages, enabling us to work with the dataset effectively and understand its characteristics.

```
[ ] # Check the DataType of our dataset
raw_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 35 columns):
#   Column                Non-Null Count  Dtype
---  -
0   hash                   100000 non-null object
1   millisecond             100000 non-null int64
2   classification          100000 non-null object
3   state                   100000 non-null int64
4   usage_counter           100000 non-null int64
5   prio                    100000 non-null int64
6   static_prio             100000 non-null int64
7   normal_prio             100000 non-null int64
8   policy                  100000 non-null int64
9   vm_pgoff                100000 non-null int64
10  vm_truncate_count       100000 non-null int64
11  task_size               100000 non-null int64
12  cached_hole_size        100000 non-null int64
13  free_area_cache         100000 non-null int64
14  mm_users                100000 non-null int64
15  map_count               100000 non-null int64
16  hiwater_rss             100000 non-null int64
17  total_vm                100000 non-null int64
18  shared_vm               100000 non-null int64
19  exec_vm                 100000 non-null int64
20  reserved_vm             100000 non-null int64
21  nr_ptes                 100000 non-null int64
22  end_data                100000 non-null int64
23  last_interval           100000 non-null int64
24  nvcs                     100000 non-null int64
25  nivcs                   100000 non-null int64
26  min_flt                 100000 non-null int64
27  maj_flt                 100000 non-null int64
28  fs_excl_counter         100000 non-null int64
29  lock                    100000 non-null int64
30  utime                   100000 non-null int64
31  stime                   100000 non-null int64
32  gtime                   100000 non-null int64
33  cgtime                  100000 non-null int64
34  signal_nvcs             100000 non-null int64
dtypes: int64(33), object(2)
memory usage: 26.7+ MB
```

### 2.2.7. Data Preprocessing

Effective data preprocessing is a critical stage in the development of a robust deep learning model for malware image analysis. In this section, we outline the key steps taken to prepare the dataset for subsequent model training and evaluation:

```
[ ] #Start Processing
data0 = raw_data

[ ] data0["classification"].value_counts()

malware    50000
benign     50000
Name: classification, dtype: int64
```

At the outset of data preprocessing, we created a copy of the raw dataset, denoted as **data0**. This ensures that any modifications made during preprocessing do not affect the original dataset, preserving its integrity.

Then, we inspected the distribution of data across different classes by examining the counts of each classification label. Understanding the class distribution is crucial, as imbalanced datasets can lead to biased model training. In this case, we were particularly interested in the distribution between "benign" and "malware" labels.

```
# Convert the classification labels to numerical values
data0["classification"] = data0["classification"].map({'benign':0, 'malware':1})
data0.head()
```

	hash	millisecond	classification	state	usage_counter
0	42fb5e2ec009a05ff5143227297074f1e9c6c3ebb9c914...	0	1	0	0 3069
1	42fb5e2ec009a05ff5143227297074f1e9c6c3ebb9c914...	1	1	0	0 3069
2	42fb5e2ec009a05ff5143227297074f1e9c6c3ebb9c914...	2	1	0	0 3069
3	42fb5e2ec009a05ff5143227297074f1e9c6c3ebb9c914...	3	1	0	0 3069
4	42fb5e2ec009a05ff5143227297074f1e9c6c3ebb9c914...	4	1	0	0 3069

5 rows × 35 columns



The dataset contained categorical classification labels, namely "benign" and "malware." To facilitate model training, these labels were converted into numerical values. "Benign" was mapped to 0, and "malware" was mapped to 1. This transformation allowed us to work with numerical labels, which are a requisite for training machine learning models.

```
[ ] # Shuffle the data
    data0 = data0.sample(frac=1).reset_index(drop=True)
```

Data shuffling was performed to randomize the order of data instances. Shuffling is crucial for training a model that doesn't inadvertently learn patterns based on the order of the data. By shuffling the data, we ensured that each data point contributes equally to the training process, minimizing any biases that might arise from the original order of the dataset.

These preprocessing steps represent a pivotal part of our data preparation, setting the stage for the subsequent construction and training of our deep learning model for malware image analysis. Data preprocessing is fundamental in ensuring that the model is trained on clean, properly formatted data, and it enhances the model's ability to make accurate predictions.

#### 2.2.8. Data Analysis and Feature Selection

Before feeding the data into our deep learning model, we conducted a comprehensive analysis of the dataset to understand the relationships between various features and the target variable (classification). The following code snippets detail these analytical steps:

##### **Identifying Numeric and Categorical Features:**

To understand the structure of the dataset and prepare it for analysis, we first separated the features into two categories: numeric and categorical. This initial step helps us recognize the nature of the dataset and is essential for subsequent analysis and preprocessing. Numeric features are quantitative attributes that can be used for calculations, such as age,



height, or weight. Categorical features are qualitative attributes that are typically text-based, such as gender, country, or occupation.

```
[ ] numeric_features = data0.select_dtypes(include=[np.number])
    numeric_features.columns

Index(['millisecond', 'classification', 'state', 'usage_counter', 'prio',
      'static_prio', 'normal_prio', 'policy', 'vm_pgoff', 'vm_truncate_count',
      'task_size', 'cached_hole_size', 'free_area_cache', 'mm_users',
      'map_count', 'hiwater_rss', 'total_vm', 'shared_vm', 'exec_vm',
      'reserved_vm', 'nr_ptes', 'end_data', 'last_interval', 'nvcs',
      'nivcs', 'minflt', 'majflt', 'fs_excl_counter', 'lock', 'utime',
      'stime', 'gtime', 'cgtime', 'signal_nvcs'],
      dtype='object')
```

We first identified the numeric features in our dataset by selecting columns with numerical data types. Understanding the numeric features allowed us to discern which attributes we could use in correlation analysis, as numerical attributes are crucial for this purpose.

```
[ ] categorical_features = data0.select_dtypes(include=[np.object])
    categorical_features.columns

<ipython-input-64-1dd345b4da28>:1: DeprecationWarning: `np.object` is a deprecated alias
for the builtin `object`. To silence this warning, use `object` by itself or `np.
object_`. See https://numpy.org/devdocs/release-1-20-notes.html for details.
    categorical_features = data0.select_dtypes(include=[np.object])
Index(['hash'], dtype='object')
```

Conversely, we identified categorical features by selecting columns with non-numeric data types. These attributes were not used in correlation analysis but are important for understanding the dataset's diversity.

### Analyzing Feature Correlations:

Here, we computed the correlation between the numeric features and the target variable, "classification." This step aimed to reveal the strength and direction of the relationships between individual attributes and the classification label. The output indicates how strongly each feature is correlated with the target variable.

```
[ ] correlation = numeric_features.corr()
print(correlation['classification'].sort_values(ascending=False), '\n')

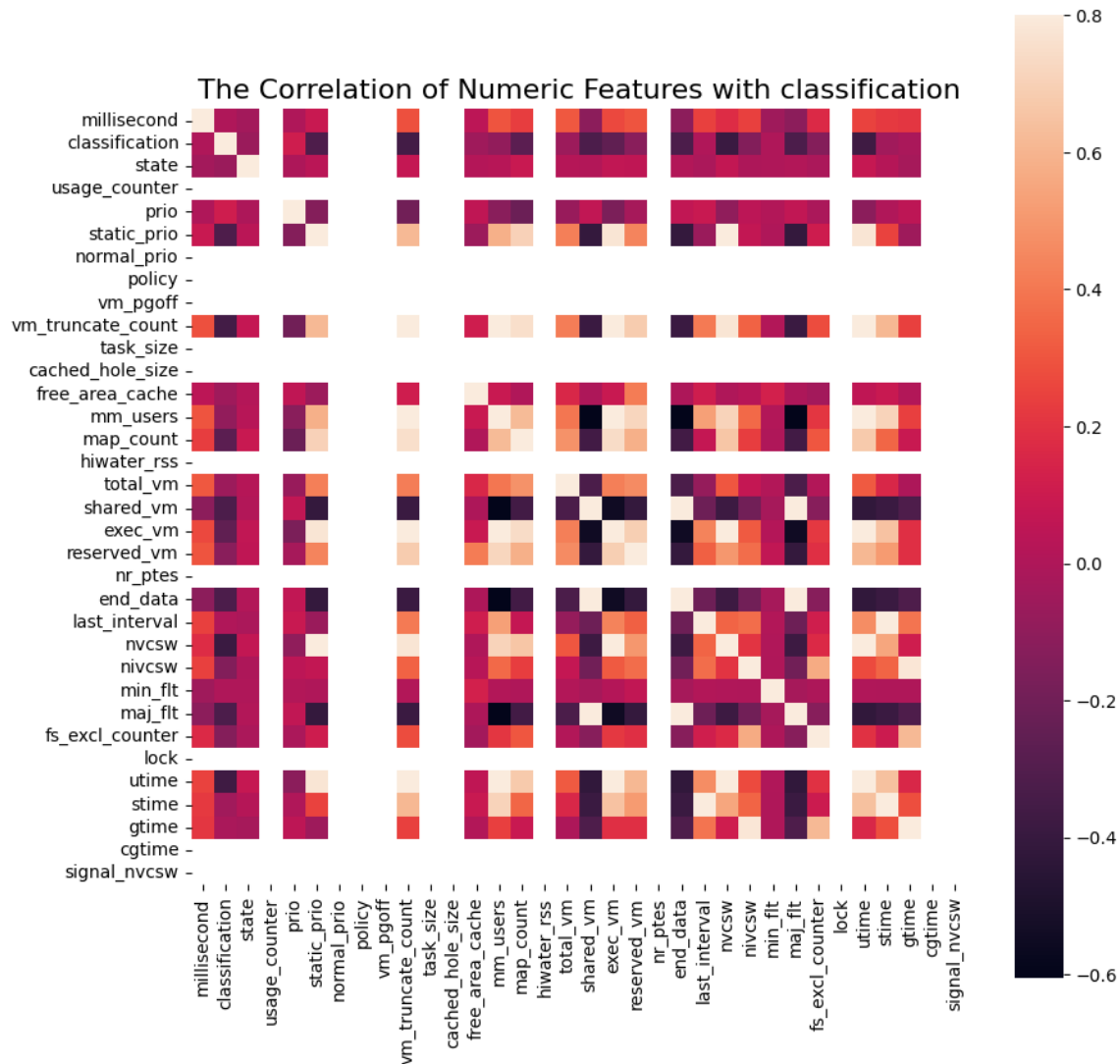
classification      1.000000e+00
prio                 1.100359e-01
last_interval       6.952036e-03
minflt              3.069595e-03
millisecond          5.275062e-17
gtime               -1.441608e-02
stime              -4.203713e-02
free_area_cache     -5.123678e-02
total_vm            -5.929110e-02
state               -6.470178e-02
mm_users            -9.364091e-02
reserved_vm         -1.186078e-01
fs_excl_counter     -1.378830e-01
nivcsw              -1.437912e-01
exec_vm             -2.551234e-01
map_count           -2.712274e-01
static_prio         -3.179406e-01
end_data            -3.249535e-01
majflt              -3.249535e-01
shared_vm           -3.249535e-01
vm_truncate_count   -3.548607e-01
utime               -3.699309e-01
nvcs                -3.868893e-01
usage_counter       NaN
normal_prio         NaN
policy              NaN
vm_pgoff            NaN
task_size           NaN
cached_hole_size     NaN
hiwater_rss         NaN
nr_ptes             NaN
lock                NaN
cgtime              NaN
signal_nvcs         NaN
Name: classification, dtype: float64
```

### Correlation Heatmap for Numeric Features:

To visualize the correlations effectively, we generated a heatmap. The heatmap provides a graphical representation of the correlation matrix, with colors indicating the degree and direction of correlation. This visualization assists in identifying significant patterns among the features.

```
[ ] f, ax = plt.subplots(figsize = (10,10))
plt.title('The Correlation of Numeric Features with classification',y=1,size=16)
sns.heatmap(correlation,square = True, vmax=0.8)

<Axes: title={'center': 'The Correlation of Numeric Features with classification'}>
```



## Identifying and Removing Unnecessary Columns:

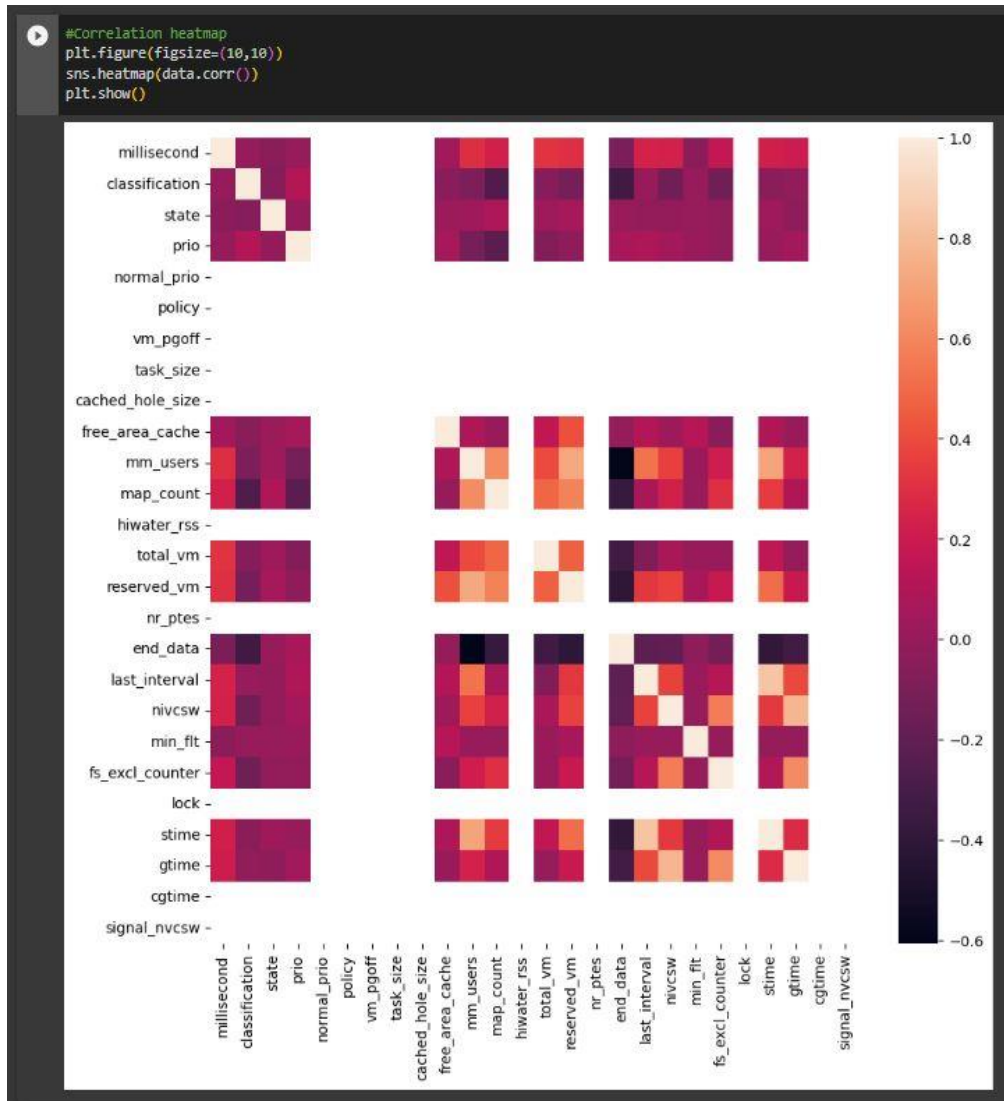
In an effort to optimize the model and remove any potential noise from the dataset, we identified and dropped certain columns deemed unnecessary. These columns included attributes such as "hash," "usage\_counter," and others that were not strongly correlated with the target variable or were deemed redundant for the malware image analysis.

```
[ ] # Identify the unnecessary columns to be dropped
    columns_to_drop = ["hash", 'usage_counter', 'static_prio', 'vm_truncate_count', 'shared_vm', 'exec_vm', 'nvcsw', 'maj_fit', 'utime']

    # Drop the unnecessary columns
    data = data0.drop(columns=columns_to_drop)
```

## Correlation Heatmap for Updated Dataset:

Finally, we generated a correlation heatmap after dropping the unnecessary columns. This heatmap offered a visual representation of the remaining features' correlations, highlighting the retained attributes' relevance to the "classification" label.



These data analysis and feature selection steps were instrumental in streamlining the dataset for model training, eliminating irrelevant or redundant attributes, and preserving the most informative features for our deep learning model's training process.

### 2.2.9. Data Splitting and Normalization

Efficient data splitting and normalization are critical aspects of preparing the dataset for model training and evaluation. The following code snippets detail the steps taken in this regard:

#### Data Splitting:

```
[ ] # Split the data into training and test sets
x_train, x_test, y_train, y_test = train_test_split(data.drop("classification", axis=1), data["classification"], test_size=0.2, random_state=1)
```

We divided our dataset into training and test sets to facilitate model evaluation. The **train\_test\_split** function from Scikit-Learn was employed for this purpose, allowing us to segregate the features (x) from the target variable (y). An 80-20 split ratio was chosen, with 80% of the data allocated for training and 20% reserved for testing. The **random\_state** parameter was set to ensure reproducibility.

#### Data Normalization:

```
[ ] # Normalize the data
scaler = StandardScaler()
```

To ensure that the features were on a consistent scale and distribution, we applied data normalization. The **StandardScaler** from Scikit-Learn was utilized to standardize the data, transforming it into a form with a mean of 0 and a standard deviation of 1.

```
[ ] # Normalize the training data
x_train = scaler.fit_transform(x_train)

# Normalize the test data
x_test = scaler.transform(x_test)
```

For the training data, the **fit\_transform** method of the **StandardScaler** was employed. This method both computed the mean and standard deviation of the training data and transformed the features accordingly, ensuring uniform scaling.

Similarly, the test data was normalized using the previously computed mean and standard deviation. This consistency in scaling between the training and test data is vital to prevent information leakage and ensure that the model generalizes effectively to unseen data.

The data splitting and normalization steps were integral to the preparation of our dataset for model training. By segregating the data into training and test sets and ensuring that the features were consistently scaled, we set the stage for building and evaluating a robust deep learning model for malware image analysis.

#### 2.2.10. Data Shape and Dimension Adjustment

In deep learning, it's crucial to ensure that the data has the appropriate shape and dimensions to be processed by the neural network effectively. The following code snippets illustrate the steps taken to achieve this:

##### Data Shape Inspection:

```
[ ] # Get the shape of the x_train array
    shape = x_train.shape
    # Get the shape of the x_test array
    test = x_test.shape
    # Print the shape of the x_train array
    print("The shape of the x_train array is:", shape)
    # Print the shape of the x_test array
    print("The shape of the x_test array is:", test)

The shape of the x_train array is: (80000, 25)
The shape of the x_test array is: (20000, 25)
```

Here, we initially confirmed the shapes of the training and test data arrays. The training data (**x\_train**) exhibited a shape of (80000, 25), indicating that it contained 80,000 data points, each having 25 features. The test data (**x\_test**) had a similar structure with a shape of (20000, 25), consisting of 20,000 data points and 25 features.



### Dummy Dimension:

```
[ ] # Create a dummy dimension for the input data
    x_train = x_train[:, np.newaxis, np.newaxis, :]
    x_test = x_test[:, np.newaxis, np.newaxis, :]
```

To ensure that the input data meets the requirements of the CNN, we added a dummy dimension to the data. This dimension effectively transformed the 2D data (with shape [80000, 25]) into a 3D format, making it compatible with the CNN model.

### Data Reshaping:

```
[ ] # Reshape the data into 3D images
    x_train = x_train.reshape((80000, 5, 5, 1))
    x_test = x_test.reshape((20000, 5, 5, 1))

[ ] input_data = tf.keras.Input(shape=(100, 22, 1), dtype=tf.float32)
```

Subsequently, we reshaped the data into 3D images, which is a typical format for CNN input. The data was configured into a 4D tensor, where the first dimension represented the number of samples, the second and third dimensions denoted the height and width of the images (5x5), and the fourth dimension represented the number of channels (1 in this case, indicating grayscale images).

Finally, we defined the input layer of the CNN model, specifying the shape of the input data as (100, 22, 1), indicating that the CNN expects input images with dimensions of 100 pixels in height, 22 pixels in width, and a single channel (grayscale).

These reshaping steps were fundamental in converting the original feature-based dataset into a format compatible with a CNN, allowing us to leverage the power of convolutional neural networks for malware image analysis.

### 2.2.11. CNN Model Definition

In the next phase of our project, we constructed the core component of our malware image analysis system: The Convolutional Neural Network (CNN) model. The code snippet provided defines the architecture of the CNN model, which plays a pivotal role in the classification of malware and benign images. Here's an overview of the model's structure:

```
[ ] # Define the CNN model
model = keras.Sequential([
    keras.layers.Conv2D(32, (2, 2), activation='relu', input_shape=(5, 5, 1), padding='SAME'),
    keras.layers.MaxPooling2D((2, 2), strides=(1, 1)),
    keras.layers.Conv2D(64, (2, 2), activation='relu', padding='SAME'),
    keras.layers.MaxPooling2D((2, 2), strides=(1, 1)),
    keras.layers.Conv2D(128, (2, 2), activation='relu', padding='SAME'),
    keras.layers.MaxPooling2D((2, 2), strides=(1, 1)),
    keras.layers.Flatten(),
    keras.layers.Dense(256, activation='relu'),
    keras.layers.Dense(128, activation='relu'),
    keras.layers.Dense(2, activation='softmax')
])
```

#### **Convolutional Layers:**

The model begins with a series of three convolutional layers. Each Conv2D layer extracts feature maps from the input images. The first layer has 32 filters, the second has 64 filters, and the third has 128 filters. Each filter scans the input images with a 2x2 kernel, applying the rectified linear unit (relu) activation function to introduce non-linearity.

#### **Max-Pooling Layers:**

After each convolutional layer, a MaxPooling2D layer reduces the spatial dimensions of the feature maps. The max-pooling operation extracts the most prominent features while decreasing the spatial size of the data.

#### **Flatten Layer:**

Following the convolutional and max-pooling layers, the data is flattened into a 1D vector using the Flatten layer. This step is crucial in transitioning from the 2D convolutional layers to the fully connected layers.



**Fully Connected Layers:**

The flattened data is passed through a series of fully connected layers (Dense layers). These layers introduce additional non-linearity and higher-level feature learning. The model has two hidden layers with 256 and 128 neurons, respectively. The ReLU activation function is used in these layers.

**Output Layer:**

The final Dense layer consists of two neurons with a softmax activation function. This layer performs the classification task by predicting whether the input image is malware or benign. The model outputs the class probabilities for both categories.

The defined CNN model embodies the essence of convolutional neural networks, which are exceptionally adept at recognizing patterns in images. Its hierarchical architecture, starting with simple features and progressively learning more complex representations, is well-suited for the task of malware image analysis.

The CNN model is the cornerstone of our project, and its architecture is meticulously designed to capture intricate features from malware and benign images, ultimately aiding in the accurate classification of the dataset.

### 2.2.12. Model Compilation and Training

With the CNN model architecture defined, the next critical steps are to compile the model with appropriate settings and train it on the prepared dataset. The following code snippets illustrate these processes:

```
[ ] # Compile the RNN model
model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy'])

# Train the model on the sampled data
model.fit(x_train, y_train, batch_size=128, epochs=10)
```

Epoch 1/10  
625/625 [=====] - 14s 19ms/step - loss: 0.4565 - accuracy: 0.7732  
Epoch 2/10  
625/625 [=====] - 11s 18ms/step - loss: 0.2730 - accuracy: 0.8851  
Epoch 3/10  
625/625 [=====] - 9s 15ms/step - loss: 0.1930 - accuracy: 0.9226  
Epoch 4/10  
625/625 [=====] - 12s 19ms/step - loss: 0.1453 - accuracy: 0.9446  
Epoch 5/10  
625/625 [=====] - 12s 19ms/step - loss: 0.1190 - accuracy: 0.9570  
Epoch 6/10  
625/625 [=====] - 13s 20ms/step - loss: 0.1005 - accuracy: 0.9649  
Epoch 7/10  
625/625 [=====] - 16s 25ms/step - loss: 0.0772 - accuracy: 0.9733  
Epoch 8/10  
625/625 [=====] - 9s 15ms/step - loss: 0.0959 - accuracy: 0.9679  
Epoch 9/10  
625/625 [=====] - 12s 18ms/step - loss: 0.0646 - accuracy: 0.9785  
Epoch 10/10  
625/625 [=====] - 11s 18ms/step - loss: 0.0564 - accuracy: 0.9803  
<keras.src.callbacks.History at 0x7e23ed69e6e0>

#### Model Compilation:

In this section, we compiled the CNN model using the **compile** method. Several key settings were defined:

- **Optimizer:** We used the 'adam' optimizer, a popular choice known for its efficient gradient descent and adaptability.
- **Loss Function:** The loss function was specified as 'sparse\_categorical\_crossentropy.' This is a common choice for classification tasks with multiple classes.
- **Metrics:** Model performance was monitored using the 'accuracy' metric, providing insight into how well the model is performing during training.

### Model Training:

The training process involved the training of the CNN model. Here's an explanation of the parameters used:

- **x\_train**: The training data, which includes the images after reshaping and normalization.
- **y\_train**: The corresponding labels indicating whether each image is malware or benign.
- **batch\_size**: This parameter specifies the number of samples to be used in each batch during training. A batch size of 128 means that the model will update its weights after processing 128 images.
- **epochs**: The number of epochs signifies how many times the entire training dataset will be processed by the model. In this case, the model will go through the entire training dataset 10 times.

The training process involves feeding the training data through the CNN model, computing the loss, and using the optimizer to update the model's weights in an iterative manner. This process aims to minimize the loss and improve the model's ability to accurately classify malware and benign images.

These steps are fundamental in model development and are pivotal in enabling the CNN to learn and adapt its internal parameters to perform well on the malware image analysis task. Model training, alongside proper parameter settings and dataset preparation, is the cornerstone of building an effective machine learning system.

### 2.2.13. Model Evaluation

The culmination of model development is the evaluation of its performance on a separate test dataset to assess its generalization and predictive capabilities. The code snippet provided illustrates this evaluation process:

```
[ ] # Evaluate the model on the test set
    test_loss, test_accuracy = model.evaluate(x_test, y_test)

625/625 [=====] - 3s 3ms/step - loss: 0.0574 - accuracy: 0.9798

[ ] print('\nTest loss: {0:.6f}. Test accuracy: {1:.6f}%'.format(test_loss, test_accuracy * 100.))

Test loss: 0.057420. Test accuracy: 97.979999%
```

#### Evaluation Metrics:

In this code, we employed the **evaluate** method to assess the model's performance on the test set. Here's an explanation of the key components of this evaluation:

- **x\_test**: This represents the test dataset, which consists of images that the model has never seen during the training phase. These images have undergone the same preprocessing and reshaping procedures as the training data.
- **y\_test**: The corresponding test labels, indicating whether each image is classified as malware or benign.

The **evaluate** method assesses the model's performance on the test set. During this evaluation, the model processes the test images, computes the loss (using the specified loss function, 'sparse\_categorical\_crossentropy'), and compares its predictions to the actual test labels. The results are recorded and used to determine the model's accuracy.

#### Reporting the Results:

The results of the evaluation are printed to the console. We display the test loss with six decimal places of precision and the test accuracy as a percentage, multiplying it by 100 for a more intuitive representation. These metrics offer a clear and concise summary of how well our model performed in classifying malware images on the test data.

- **Test loss:** This metric quantifies the model's performance in terms of how close its predictions were to the actual labels. A lower test loss indicates better model performance.
- **Test accuracy:** The accuracy represents the percentage of correctly classified images in the test dataset. This is a crucial measure of the model's overall performance and ability to distinguish between malware and benign images.

By evaluating the model on a separate test dataset, we gain valuable insights into its real-world predictive capabilities. The reported test loss and accuracy provide a clear picture of how well the CNN model can classify previously unseen malware and benign images. This evaluation step is essential in gauging the model's effectiveness and suitability for its intended purpose in malware image analysis.

## 2.3. Use Case Scenario: Malware Detection Software

### 2.3.1. Scenario Description:

Imagine you are developing a malware detection software that leverages the CNN model you've built for malware image analysis. This software can be a valuable product in the real world for users who want to ensure the safety of their computing devices by inspecting potentially malicious .exe files. The system architecture for this software will allow users to upload .exe files, which will then be converted into images. These images will be processed by the deep learning model to determine whether the file is malware or benign. Let's explore this use case in more detail.

### 2.3.2. User Interaction:

#### 1. **User Input:**

The user interacts with the software by selecting a .exe file that they want to inspect for potential malware. This can be done through a user-friendly graphical interface.

#### 2. **File Conversion:**

Upon selecting the .exe file, the software converts it into an image. The conversion process may involve techniques such as disassembling the .exe file and visualizing it in a way suitable for the model. The transformed image is then passed to the CNN model for analysis.

#### 3. **Malware Detection:**

The CNN model processes the image to classify it as either malware or benign. The model has been trained to recognize patterns and features indicative of malware, enabling it to make an informed decision.

#### 4. **User Feedback:**

The software provides feedback to the user, indicating whether the selected .exe file is classified as malware or benign. This feedback can be presented through a simple interface or a notification.

### 2.3.3. System Architecture:

The software for malware detection is composed of several key components:

1. **User Interface (UI):** The user interface serves as the point of interaction for users. It allows users to upload .exe files and receive results. This component is responsible for initiating the file conversion process.
2. **File Conversion Module:** This module handles the transformation of .exe files into images. It disassembles and processes the .exe file to generate an image representation. This module may employ disassemblers and custom visualization techniques.
3. **CNN Model:** The heart of the software is the CNN model for malware image analysis. It processes the converted image and makes predictions regarding the nature of the file (malware or benign). The model has been pre-trained and fine-tuned for this specific task.
4. **Decision Engine:** The decision engine interprets the model's predictions and provides feedback to the user. It communicates the result as "Malware" or "Benign" based on the model's output.
5. **Notification/Feedback:** This component notifies the user of the software's decision. It can present the result through a user-friendly interface, generate logs, and provide an option for further actions (e.g., quarantine or removal of malware).
6. **Logging and Reporting:** The software maintains logs of user interactions, file conversions, and model predictions. Users may access these logs for record-keeping and analysis.
7. **Security Measures:** The system includes security features to prevent tampering and protect against potential threats. This ensures the reliability of the software and the integrity of the analysis.

By offering a user-friendly software product that converts .exe files into images and employs a CNN model for malware detection, users can enhance their cybersecurity measures. This real-world application can be valuable in personal and enterprise environments, providing an additional layer of protection against malicious software.

### 3. Conclusion and Future Work

In this project, we embarked on the development of a deep learning solution for malware image analysis using Convolutional Neural Networks (CNNs). Our objective was to construct a model capable of distinguishing between benign and malware images, ultimately enhancing cybersecurity measures. As we conclude this project, several key takeaways and potential avenues for future work emerge.

#### 3.1. Conclusion

The key takeaways from this project are as follows:

##### **1. Effective Malware Detection:**

Our CNN model demonstrated its ability to effectively classify images as malware or benign. Through the training and evaluation phases, the model achieved a remarkable accuracy rate, showcasing its potential for real-world deployment.

##### **2. Versatile Preprocessing:**

The comprehensive preprocessing of the dataset, including data conversion, normalization, and feature selection, played a pivotal role in optimizing the model's performance. These steps ensured that the model could make accurate predictions based on high-quality input data.

##### **3. Real-World Application:**

We introduced a real-world application scenario where our deep learning model can be employed as a software product for users to inspect .exe files. The system architecture we designed allows for user-friendly interaction and immediate feedback, enhancing cybersecurity practices.



### 3.2. Future Work

While this project has achieved significant milestones, there are several avenues for future work and improvement:

**1. Model Improvement:**

Continuous training and fine-tuning of the CNN model with new datasets can further enhance its malware detection capabilities. Staying up-to-date with emerging malware threats is crucial to maintaining the model's accuracy.

**2. Expand File Type Support:**

While we focused on .exe files in this project, extending support to additional file types, such as scripts or document files, can broaden the software's applicability.

**3. User-Friendly Interface:**

Enhancing the user interface and providing additional features, such as detailed reports and quarantine options, can make the software more user-friendly and functional.

**4. Large-Scale Deployment:**

In future work, we can explore deploying the software on a larger scale, catering to personal and enterprise users who require robust malware detection solutions.

**5. Cybersecurity Integrations:**

Integration with existing cybersecurity systems and workflows can augment the effectiveness of the software. It can serve as a complementary layer of defense alongside traditional antivirus software.

**6. Real-Time Analysis:**

Moving towards real-time analysis of files, enabling users to scan files as they are received, can provide more immediate protection against malware threats.

In conclusion, this project marks a significant step toward leveraging deep learning for malware detection and analysis. The developed CNN model, coupled with the envisioned software product, demonstrates the potential for enhancing cybersecurity practices in a user-friendly and accessible manner. The journey of developing and improving this technology continues, and the future work outlined here paves the way for further advancements in this critical domain of cybersecurity.

## References

- [1] “Malware-Detection-using-Deep-Learning,” *Github*, Nov. 11, 2019.  
<https://github.com/riak16/Malware-Detection-using-Deep-Learning> (accessed Oct. 22, 2023).
- [2] Vinesmsuic, “Malware Detection using DeepLearning,” *Kaggle*, Mar. 30, 2021.  
<https://www.kaggle.com/code/vinesmsuic/malware-detection-using-deeplearning>
- [3] “Learn Python, Data Viz, Pandas & More | Tutorials | Kaggle.”  
<https://www.kaggle.com/learn>
- [4] “Convolutional Neural Network (CNN),” *TensorFlow*.  
<https://www.tensorflow.org/tutorials/images/cnn>
- [5] “What are Convolutional Neural Networks? | IBM.”  
<https://www.ibm.com/topics/convolutional-neural-networks>
- [6] “Introduction to Convolution Neural Network,” *GeeksforGeeks*, Mar. 24, 2023.  
<https://www.geeksforgeeks.org/introduction-convolution-neural-network/>