# VIRUS DETECTION USING DATA MINING

Submitted in partial fulfilment

Of

**Mini project** in Bachelor of Technology

Submitted by

1. **JAYA SHANKAR REDDY KOMMURU 411529**

2. **VISHNU NAGINENI 411551**

3. **SAI SHIVA KUMAR KATHRI 411535**

4. **V. PRATHYUSH 411576**

Under the Supervision of

**Mr. BKSP Kumar Raju Alluri**



# NATIONAL INSTITUTE OF TECHNOLOGY,

# Andhra Pradesh.

### SUMMARY

❖ **The project consists of two types of modules. One is the Implemented modules and other is the Theoretical Modules. There are 6 Implemented modules and 3 Theoretical Modules.**

❖ **The Implementation Modules use two different features for classification of malware and non-malware files, namely hex dump and System call sequences.**

❖ **The hex dump feature is used in Static Analysis where different classification algorithms of data mining are applied and their accuracy is being compared.**

❖ **Once identifying that the virus is present, the same hex dump features are used to identify which type of virus has injected that program.**

❖ **The System call sequences feature is used in Dynamic Analysis where clustering Algorithms are applied for the non-malware data set. And any outlier detected is treated as a Malware.**

❖ **Instead of considering all the sequences, the top k malware sequences are found using the TKS Algorithm. Now these top k sequences are compared with a given sequence and found out that whether it is malicious or genuine.**

❖ **Considering an assumption that the malware sequences are complete, from a given sequence, all the in between calls are counted apart from the malware calls. The trace having huge count is more probable to be malicious.**

❖ **Considering another assumption that there are fixed number of genuine activities for an Operating System, all the in between calls of a sequence are appended and checked for genuine activities. More the genuity, more the probability of it being non malicious.**

❖ **Theoretical model of integrating the static and dynamic model is given.**

❖ **Two different tools are explored for dynamic Analysis, namely Cuckoo sandbox and Intel's PIN tool.**

❖ **PIN tool is used for detecting ROP.**

❖ **Cuckoo sandbox is used for extracting run time features other than System call Sequences.**

## DATA MINING:

**Data mining** is the computing process of discovering patterns in large data sets involving methods at the intersection of machine learning, statistics, and database systems. It is an essential process where intelligent methods are applied to extract data patterns. It is an interdisciplinary subfield of computer science. The overall goal of the data mining process is to extract information from a data set and transform it into an understandable structure for further use. Aside from the raw analysis step, it involves database and data management aspects, data pre-processing, model and inference considerations, interestingness metrics, complexity considerations, post-processing of discovered structures, visualization, and online updating. Data mining is the analysis step of the "knowledge discovery in databases" process, or KDD.

The actual data mining task is the semi-automatic or automatic analysis of large quantities of data to extract previously unknown, interesting patterns such as groups of data records (CLUSTER ANALYSIS), unusual records (ANAMOLY DETECTION), and dependencies (ASSOCIATION RULE MINING,SEQUENTIAL PATTERN MINING). This usually involves using database techniques such as SPATIAL INDICES. These patterns can then be seen as a kind of summary of the input data, and may be used in further analysis or, for example, in machine learning and predictive analytics. For example, the data mining step might identify multiple groups in the data, which can then be used to obtain more accurate prediction results by a DECISION SUPPORT SYSTEM. Neither the data collection, data preparation, nor result interpretation and reporting is part of the data mining step, but do belong to the overall KDD process as additional steps.

## METHODS IN DATA MINING:

1. Preprocessing techniques.
2. Association Mining.
3. Classification Mining.
4. Cluster analysis.
5. Sequence Mining.

## MALWARE

The term malware is used when referring to malicious software, which is designed to disrupt or gain unauthorized access to a system. There are various types of malware, out of which our project concentrates on virus.

**Virus**: It is a malicious code which gets attached to normal files in a computer and remains idle until the files are opened/executed. Once executed, the virus gets multiplied and attaches to other non-infected files and destroys the functionality of the computer.

Hence, it is very essential to differentiate between a malware and a cleanware once an unknown file has entered into the system inorder to ensure it will not perform any malicious activity.

This project aims to perform malware detection, where models are trained on different kinds of data sets. These models are trained using different data mining algorithms and tested using cross validation and percentage-split.

## METHODS USED IN PROJECT:

### ❖ Preprocessing techniques:

- Hex dump generation.
- String to word vector.
- n grams tokeniser.

### ❖ Classification Mining:

- Naive Bayes**:** The Bayesian Classification represents a supervised learning method as well as a statistical method for classification. Assumes an underlying probabilistic model and it allows us to capture uncertainty about the model in a principled way by determining probabilities of the outcomes. It can solve diagnostic and predictive problems.
- Naive Bayes Multinomial: The multinomial Naive Bayes classifier is suitable for classification with discrete features (e.g., word counts for text classification). The multinomial distribution normally requires integer feature counts.
- NaiveBayesMultinomial Text: Multinomial Naive Bayes is a specialized version of Naive Bayes that is designed more for text documents. Whereas simple naive Bayes would model a document as the presence and absence of particular words, multinomial naive bayes explicitly models the word counts and adjusts the underlying calculations to deal with in.
- NaivebayesMultinomial updatable : This version of Naïve Bayes takes the input increamentally and keeps updating the model for each tuple.
- Naive Bayes Updateable :This version is similar to the above one.
- J-48: It is java implementation of C 4.5 in weka.
- Random Forest: Random Forests are a combination of tree predictors where each tree depends on the values of a random vector sampled independently with the same distribution for all trees in the forest. The basic principle is that a group of "weak learners" can come together to form a "strong learner".
- Random Tree: This version comes under the Random Forest Algorithm.
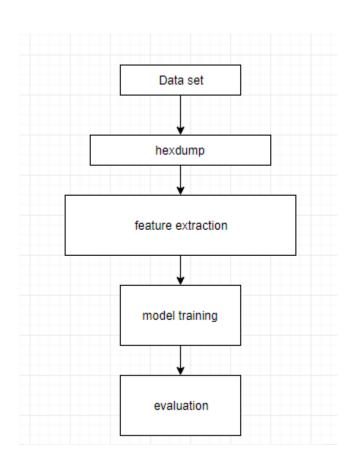
### ❖ Clustering:

- Agglomerative Hierarchical Clustering (single, complete, average, mean, center, median): Agglomerative Hierarchical clustering -This algorithm works by grouping the data one by one on the basis of the nearest distance measure of all the pairwise distance between the data point.
- Simple K means: Forms K clusters. All the tuples in a cluster have same behavior. K is the input.
- Density Based Clustering: Density based clustering algorithm has played a vital role in finding non-linear shapes structure based on the density. Density-Based Spatial Clustering of Applications with Noise (DBSCAN) is most widely used density based algorithm. It uses the concept of density reachability and density connectivity.

❖ **Sequence Mining:**
- TKS Algorithm: TKS is an algorithm for discovering the top-k most frequent sequential patterns in a sequence database. The input of TKS is a sequence database and a user-specified parameter named $k$ (a positive integer representing the desired number of patterns to be found).

## IMPLEMENTED MODULES:

## MODULE 1: STATIC ANALYSIS

**Static Analysis**: Without executing the malware, many features like strings embedded in binary,PE header, disassembly can be extracted from a file using debugging tools. The dataset which we considered contains directly the hex dump.

**Data set**: https://www.kaggle.com/c/malware-classification/data

**Description**: This data set contains only malware files. The files consist of **hex dump** of different malwares of windows.

Non-malware files are also needed. So, around 100 windows executable files are taken and hex dump of them is generated using a C++ program.

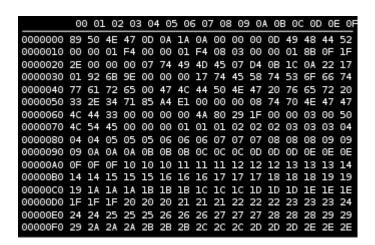**Feature used**: Hex dump is the hexadecimal representation of the instructions of an executable.



Fig1: Hex dump of a file. The starting 4 bytes of each line is the memory address followed by the hex dump.

## N Gram:

An N-gram is a sequence of bytes of fixed or variable length, extracted from the hexadecimal dump of an executable program. Here, we are using n-grams as a general term for both overlapping and non-overlapping byte sequences.

**EXAMPLE:**

**Hex Dump:** 89 50 4E 47 0D 0A 1A

If N=2 (known as bigrams), then the n grams would be:

89 50

50 4E

4E 47

47 0D

0D 0A

A 1A

If N=3 (known as trigrams), then the n grams would be :

89 50 4E

50 4E 47

4E 47 0D

47 0D 0A

0D 0A 1A

If N=4(known as four grams), then the n grams would be :

89 50 4E 47

50 4E 47 0D

4E 47 0D 0A

47 0D 0A 1A

**Text Directory Loader**

Loads all text files in a directory and uses the subdirectory names as class labels. The content of the text files will be stored in a String attribute; the filename can be stored as well.

Example directory layout:

```
...
 |
 +- text_example
    |
    +- class1
    |  |
    |  + file1.txt
    |  |
    |  + file2.txt
    |  |
    |  ...
```

```
    |
 +- class2
 |   |
 |   + another_file1.txt
 |   |
 |   + another_file2.txt
 |   |
 |   ...
```

The above directory structure can be turned into an ARFF file like this:

**weka.core.converters.TextDirectoryLoader:**

Loads all text files in a directory and uses the subdirectory names as class labels. The content of the text files will be stored in a String attribute, the filename can be stored as well.

**weka.filters.unsupervised.attribute.StringToWordVector:**

Converts String attributes into a set of attributes representing word occurrence (depending on the tokenizer) information from the text contained in the strings. The set of words (attributes) is determined by the first batch filtered (typically training data). Hex dump generation using command in Weka.

**Pre-processing steps**:

1. All the malwares' hex dumps are placed in a folder named malware. Similarly all non-malwares' hex dumps are placed in a folder named non malware.
2. Both the above folders are placed in a parent folder.
3. The parent folder is loaded into Weka using simple CLI. As simple CLI is a bit faster than GUI.
4. The arff file generated is loaded into Weka. A filter in unsupervised attribute String to Word Vector is used. Again in this a tokeniser called n gram tokeniser is used.
5. Once the above filter is applied, all the hex dumps are converted into n grams and the data set converts to a binary data set.

**Training the model:**

1. Once the data set being pre-processed, the classification algorithms are applied.

**Testing the model:**

1. The data set is tested using both percentage split and cross validation.
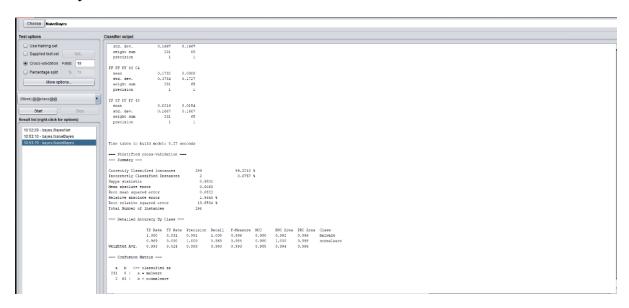
**Measurement metric:**

1. The metric used for measuring the accuracy is the confusion matrix for each algorithm.
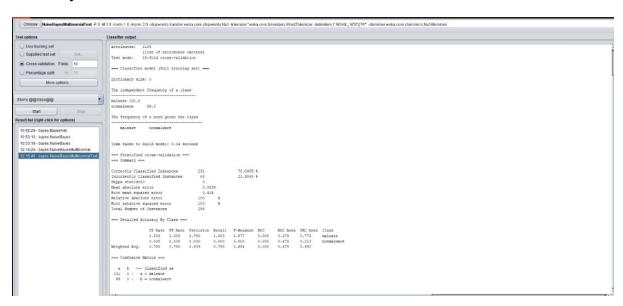
**Evaluation:**

**5-Grams:**

**CROSS-VALIDATION: 10 Folds**

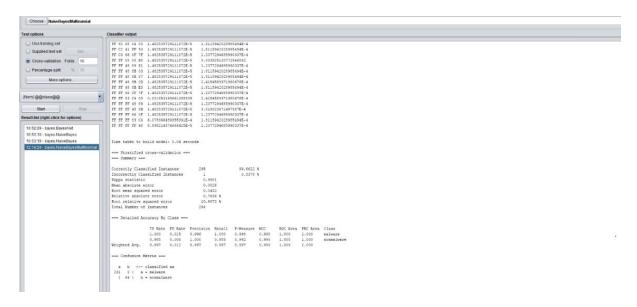**Naïve-Bayes: 99.3243%**



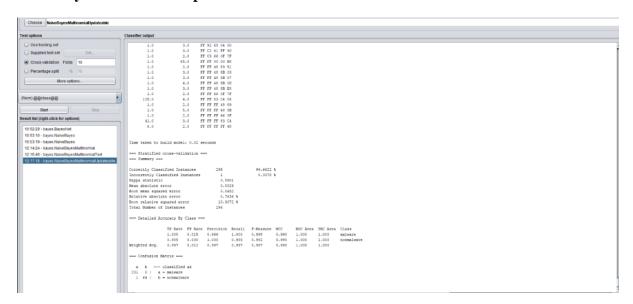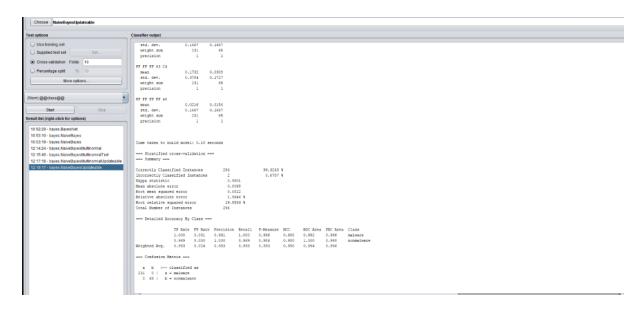**NaiveBayesMultinomialText:78.0405%**
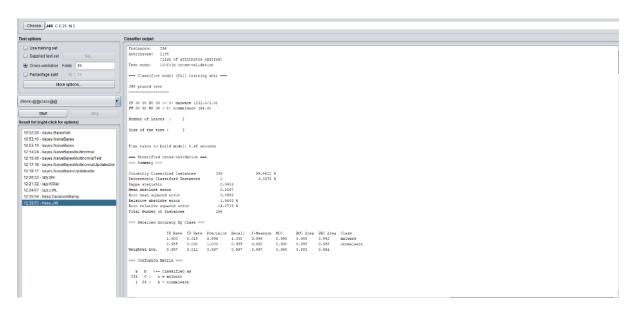


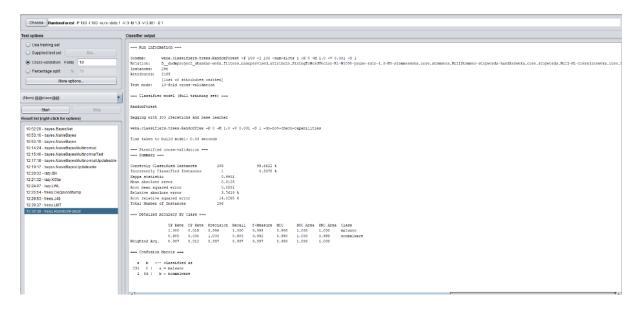**NaiveBayesMultinomial: 99.632%**

**NaiveBayesMultinomialUpdatable: 99.632%**



**NaiveBayesUpdatable:99.3243**
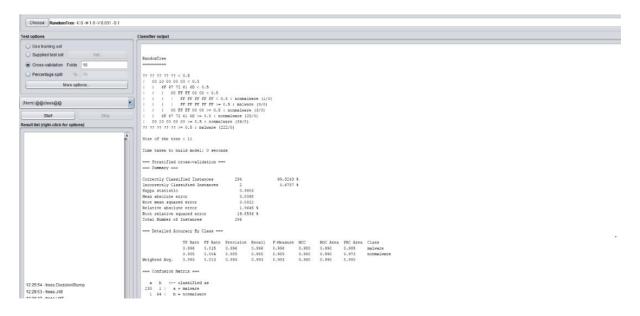
**J48: 99.6632%**



**Random Forest: 99.662%**

**RandomTree: 99.3243**



## MODULE 2: (STATIC ANALYSIS TO FIND THE TYPE OF VIRUS)

**Data set**: The Data set which was considered above contains all malware files along with the type of virus that the file has got injected with it. In the data set there are totally nine classes of viruses, namely:

1. Ramnit
2. Lollipop
3. Kelihos_ver3

4. Vundo
5. Simda
6. Tracur
7. Kelihos_ver1
8. Obfuscator.ACY
9. Gatak

**Pre-processing steps:**

1. The pre-processing steps are the same as for the above module, i.e generating n grams using String to word vector.

**Training the model:**

1. Once the data set being pre-processed, the classification algorithms are applied.

**Evaluating:**

1. Once the model being trained, the hex dump which was found to contain virus in module 1 is applied to this trained model.
2. The model, then assigns a class to this hex dump i.e. the type of virus that has injected to it.

## MODULE 3: (DYNAMIC ANALYSIS)

**Dynamic Analysis**: The features are extracted from executable during execution.

**Data set**: http://www.cs.unm.edu/~immsec/data/live-sendmail.html

**Features used**: System call sequences.

**SYSTEM CALL SEQUENCES**:

An application programming interface (API) is a source code interface that an operating system or library provides to support requests for services to be made of it by computer programs. For operating system requests, API is interchangeably termed as system calls. An API call sequences captures the activity of a program and, hence, is an excellent candidate for mining of any malicious behaviour.

Data set description: The data set consist of both malware and non-malware

Process' system call sequences generated using Linux commands like send mail, login, etc., The system calls are represented as numbers and hence the sequences as string of numbers.

**Pre-processing:**

1. All the non-malwares process' system call sequences are stored in a file named malware.
2. Generate nonmalware.arff using texttodirectory loader.
3. Similarly generate another malware. Arff for malware process' system call sequences.

**Training the model:**

1. Load the arff into Weka and apply cluster analysis.

**Testing the model**:

1. Supply the malware. Arff as the test set and run the model.
2. Detect the outliers as the malwares' system call sequences will be different from the non-malwares' system call sequences.

**Measurement metric**:

1. More the ratio of number of outliers and total number of malware files, more is the accuracy of the trained model.

# MODULE 4: (DYNAMIC ANALYSIS)

**Assumption:** The malware sequence is assumed to be of fixed length.
**Algorithm:**

1. Taken the malware system call sequences and converted to spmf format.

2. Applied the TKS Algorithm and obtained the top malware sequences.

3. Now generated the same length sequences as that of the topk sequences.

4. Each such sequence is compared with all the top k sequences. If match is found then the count of the process is incremented.

5. More the count, more the probability of being it a malware process.

# MODULE 5: (DYNAMIC ANALYSIS)

**Assumption:** The Malware system call sequences are complete and the attacker does intrusion without wasting time.

**Algorithm:**

1. Taken a process' system call sequence 's'.
2. For each Malware system call sequence $s_i$ , check whether all calls are present or not in 's'
3. If yes, check whether the same sequence of calls is present or not.

4. Else, check with other Malware system call sequence until no malware sequence available.
5. If yes, count the number of calls in between the malwares' calls.
6. Else, 's' is non-malware.
7. More the count, less the probability of it being a malware.

## MODULE 6: (DYNAMIC ANALYSIS)

**Assumption**: The Malware sequences are complete and each genuine activity's sequences are complete.

**Algorithm**:

1. Taken a process' system call sequence 's'.
2. For each Malware system call sequence $s_i$, check whether all calls are present or not in 's'
3. If yes, check whether the same sequence of calls is present or not.
4. Else, check with other Malware system call sequence until no malware sequence available.
5. If yes, append all the calls (say new sequence 'u') in between the malwares' calls.
6. Else, 's' is non-malware.
7. For each activity '$a_i$', generate n grams of 'u' where 'n' is the length of each activity.
8. Check whether each gram is equal to the '$a_i$'.
9. If yes increase the count.
10. Check with other activities until there are no activities left.
11. More the count, less the probability.

## THEORETICAL MODULES:

## MODULE 1: (INTEGRATION OF STATIC AND DYNAMIC ANALYSIS)

**Explanation:** Both Static and Dynamic Analysis techniques are complementary to each other and.Static Analysis does the complete code coverage of a program whereas Dynamic Analysis does less code coverage. But Static Analysis won't be able to provide sufficient information compared to Dynamic Analysis. They have their own advantages and disadvantages. This led to the belief that integrating both the techniques would lead us to better accuracy when compared to implementing them individually.

For an executable, if it contains malicious code in it, then it would be waste of time to perform Dynamic Analysis ( Executing the program in Virtual Machine, extracting the features and then comparing to the model trained). Instead Static Analysis can be applied and it can be detected as malicious or not.

So, first Static Analysis can be applied and then the Dynamic Analysis.

**Steps:**

1. Given a sample executable to test.
2. First Static Analysis is applied.
3. If found to be virus, then done.
4. Else, apply Dynamic Analysis.
5. Extract features and test with trained model.
6. If found to be malicious, done.
7. If not, the executable is free from viruses.

# MODULE 2: (USING INTEL'S PIN TOOL TO DETECT BASIC ROP Attacks)

**PINTOOL:** It is API to extract features from an executable during its execution.

Programs are written on top of this API to extract the desirable features**.**

## How ROP Attack came into picture:

**Attack 1:** Generally, the formal way of attacking is to inject malicious code(,i.e Payload) into any program's address space and perform malicious activity using that injected code. (This is achieved through Buffer Overflow Attack).

**Response**: To overcome this, many operating systems provide the NX bit enabled, i.e, the programs' address space is either executable or writable but not both. So, even if the buffer is overwritten with payload, it will not be executed with this NX bit enabled.

**Attack 2:** With this attackers, had only one way to perform malicious activity, i.e. using the instructions present in the DLL ( Dynamic Linked Library).But using those instructions, the extent of intrusion that can be done was limited.

**Response**: ASLR (Address Space Layout Randomisation). The addresses of the libraries were randomised. And the attackers had to apply brute force approach to find the instructions of the library.

**Attack 3:** Now comes ROP, which just uses small instructions sets already available in a program. Those small instructions are named gadgets for ROP. The Gadgets contain mostly 3-5 instructions and would end with a return instruction and hence the name Return Oriented

Programming (ROP). Now all those Gadgets are arranged in such a way that upon execution of sequence of Gadgets would lead to a malicious activity. In this way the attackers can perform any malicious activity as desired.

## Our Approach to detect ROP:

Reports and how ROP attacks programs show that the number of return instructions executed in an ROP attack and normal program execution vary greatly.

## Steps:

1. Perform normal program execution by giving normal inputs to the program and find out the number of return instructions using a program written on top of PINTOOL API.
2. Now perform execution of the program giving ROP Payload and find out the number of return instructions.
3. If the count is much greater than that of the normal executions. Count the number of instructions in between the return instructions considering a sliding window using another program written on top of PINTOOL.
4. If there is a sudden decrease in the number of instructions and the decrease continuous for more than a threshold, then it is more probable that the given input is a ROP payload and performs a malicious activity.
5. Else, it is less probable to be ROP payload.

## MODULE 3 (DYNAMIC ANALYSIS USING CUCKOO SANDBOX)

Different run time features can be extracted from the reports generated by Cuckoo Sandbox after a sample is given to the sandbox. The sandbox executes the program in a virtual Machine and submits reports to the host machine in a JSON format.

Each Analysis report contains the execution trace of a program which includes process tree, file and memory changes, memory dump and mutexes, etc.,

From the JSON format, different kinds of run time features can be extracted and can be used for training our models.

## NOVELTY:

❖ If module 1 or 2's implementation of the Theoretical models explained above give accuracy, then it leads to Novelty.

# CONCLUSION

Virus detection has been done using two different features with two different data sets. The accuracies of the different techniques applied are compared for the two features. Explored different techniques how an attack is done and different tools for Dynamic Analysis.

# FUTURE WORK

- ❖ All the Theoretical Models explained above will be implemented.
- ❖ We will be building a model to detect ROP with better accuracy.
- ❖ Also after detecting the virus and virus' type, we will try to delete that virus and ensure normal execution to the user.
- ❖ If everything above are done, will build a anti-virus software for windows OS.

## TIMELINE

| NO OF DAYS SPENT | TASK |
|---|---|
| 2 | Finding, downloading and extracting the data set for hex dump |
| 7 | Pre-processing |
| 7 | Training and Testing (module 1 completed) |
| 7 | PINTOOL AND ROP |
| 1 | Finding and Downloading the dataset for system call sequences |
| 2 | Module 2 completed |
| 4 | Module 3 completed |
| 3 | Module 4 completed |
| 2 | Module 5 completed |
| 3 | Module 6 completed |
| 3 | Cuckoo Sandbox installation and setting up environment |
| 1 | For Novelty |

**Total Number of days : 42 (Approx)**

**REFERENCES:**

1. Base paper 1:
   http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.471.7848&rep=rep1&type=pdf
2. Base paper 2:
   http://stars.library.ucf.edu/cgi/viewcontent.cgi?article=4709&context=etd
3. ROP : http://codearcana.com/posts/2013/05/28/introduction-to-return-oriented-programming-rop.html
4. PIN TOOL: https://software.intel.com/en-us/articles/pin-a-dynamic-binary-instrumentation-tool
5. CUCKOO SANDBOX: https://www.cuckoosandbox.org/

**CODES:**

**MODULE 5:**

**Java Code :**

```java
import java.io.*;

import java.util.*;

public class Test

{

public static void main(String args[])

{

String content = new Scanner(new File("filename")).useDelimiter("\\Z").next();

String content2 = new Scanner(new File("filename")).useDelimiter("\\Z").next();

String[] words = content.split("\\s+");

String[] words2= content2.split("\\s+");

int found;

for(int i=0;i<words.length;i++)

{

found =0;

  for(int j=0;j<words2.length;j++)

  {

  if(words[i].equals(words2[j]))

   {

   found=1;
```

```java
    }
  }
if(found ==0)
break;


}
if(found ==0)
{
System.out.println("all the system call sequences are not present");}
else
{
prev = -1
cur = 0;
for(int i=0;i<words.length;i++)
{
  for(int j=0;j<words2.length;j++)
   {
     if(words[i].equals(words2[j]))
      { cur = j;break;}}
if(cur<prev)
{found =0;break;
}
prev=cur;}
if(found==0)
{
System.out.println("all the calls are presentn but the sequence is missing");
}
```

```java
else

{

m=0;

count =0;

for(int i=0;i<words.length;i++)

{

while(!(words2[m].equals(words[i])))

{

++count;


++m;

}

++m;

}

}

}}
```

**MODULE 6:**

**Java Code :**

```java
import java.io.*;

import java.util.*;

public class Test2

{
```

```java
public static void main(String args[])
{
String content = new Scanner(new File("filename")).useDelimiter("\\Z").next();
String content2 = new Scanner(new File("filename")).useDelimiter("\\Z").next();
String[] words = content.split("\\s+");
String[] words2= content2.split("\\s+");
int found;
for(int i=0;i<words.length;i++)
{
found =0;
  for(int j=0;j<words2.length;j++)
  {
  if(words[i].equals(words2[j]))
   {
   found=1;
    }
  }
if(found ==0)
break;
}
if(found ==0)
{
System.out.println("all the system call sequences are not present");
}
else
{
prev = -1
```

```java
cur = 0;
for(int i=0;i<words.length;i++)
{
  for(int j=0;j<words2.length;j++)
   {
     if(words[i].equals(words2[j]))
      {
      cur = j;break;  }}
if(cur<prev)
{found =0;break;
}
prev=cur;}
if(found==0)
{
System.out.println("all the calls are presentn but the sequence is missing");
}
else
{
m=0;
count =0;
for(int i=0;i<words.length;i++)
{
StringBuffer sb = "";
while(!(words2[m].equals(words[i])))
{


sb.append(words[m]);
```

```
            ++m;
        }
        if(sb.length()!=0)
        {
        System.out.println(sb);
        }
        ++m;
    }
}}}
```

```
        if(sb.length()!=0)
```