# NATIONAL INSTITUTE OF TECHNOLOGY KARNATAKA
## Department of Computer Science and Engineering



# ARTIFICIAL INTELLIGENCE IN ANTI-VIRUS TECHNOLOGY

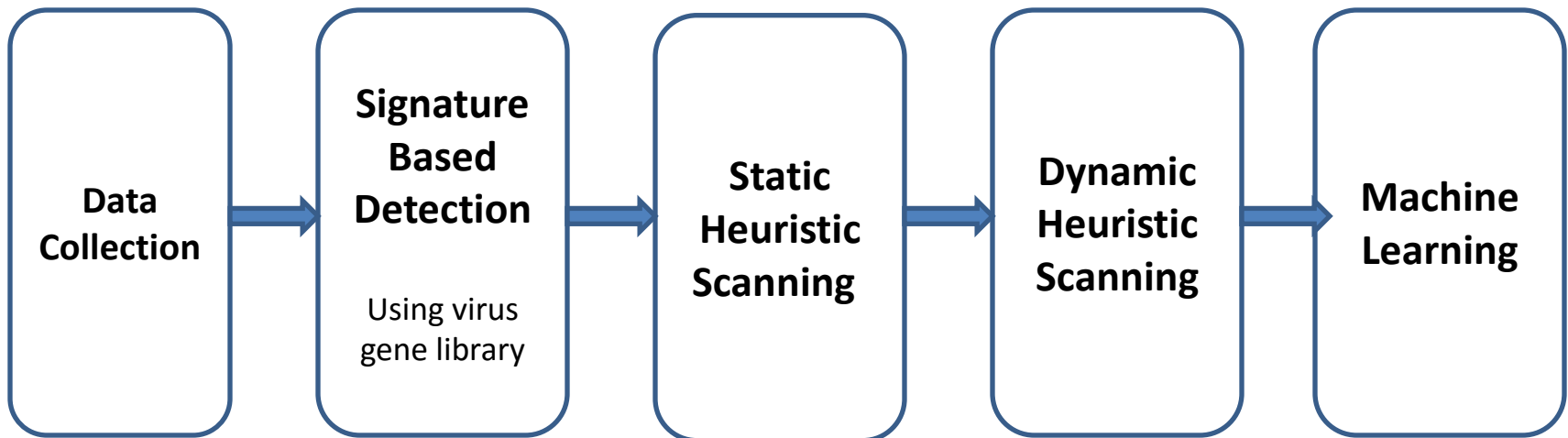**Guide :     Mr K Vinay Kumar**

**By:        Team 09UG19**
**Madhuri Shanbhogue     (09CO47)**
**Adarsh Pradhan JMT     (09CO04)**
**M C Harshavardhana     (09CO45)**

# PROBLEM DEFINITION AND OBJECTIVE

- Problem Definition: Currently used AI techniques detect large number of false positives

- Objective: Reduce the number of false positives using Learning techniques.

| Data Collection | → | **Signature Based Detection**<br><br>Using virus gene library | → | **Static Heuristic Scanning** | → | **Dynamic Heuristic Scanning** | → | **Machine Learning** |
|---|---|---|---|---|---|---|---|---|

# SIGNATURE BASED DETECTION

- Use of virus gene library

- 3,00,000 signatures

- Advantages
  - Fast
  - Available signatures

- Disadvantages
  - Cannot detect viruses without signatures

# STATIC HEURISTIC SCANNING

- Search for behavior by static code analysis
- Challenges
    - Metamorphic ( Obfuscation )
    - Polymorphic ( Encryption )
- Obfuscation types
    - Register reassignment
    - Dead Code Insertion
    - Code Transposition
    - Code Substitution
- Obfuscation Handling
    - Regular Expression
    - Control Flow Graph
- Advantages
    - Not computation intensive

```
Code snippet
E800 0000 00(90)*
5B(90)* 8D4B
42(90)* 51(90)*
50(90)* 50(90)*
0F01 4C24 FE(90)*
5B(90)* 83C3 C(90)*
FA(90)*8B2B
```
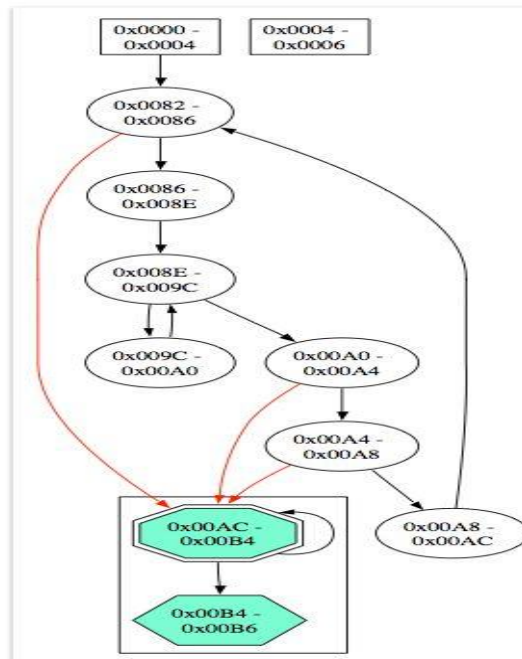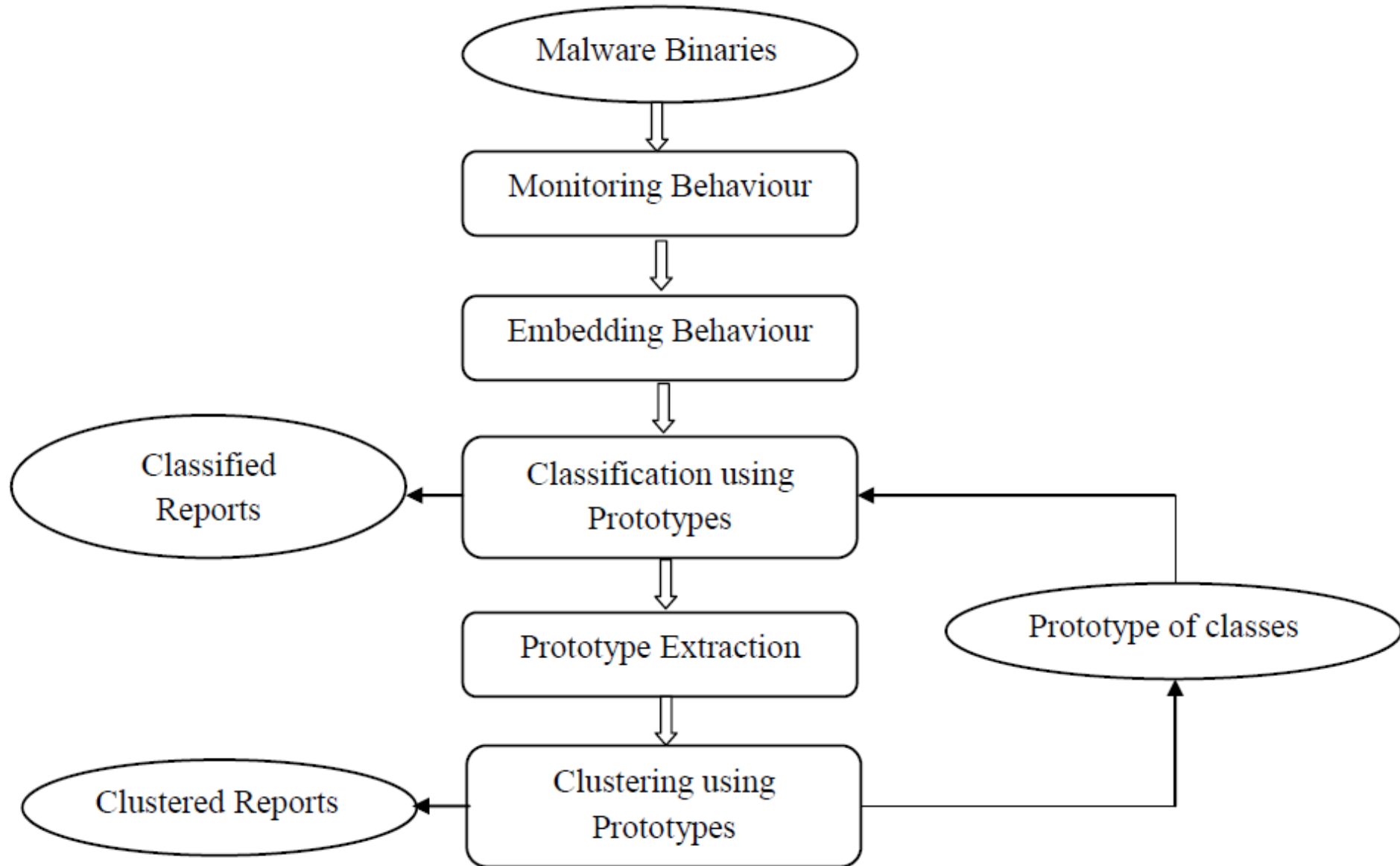
# DYNAMIC HEURISTIC SCANNING

- Decision parameters
  - System calls
  - Unusual destination
  - Analysis of file types and file system
  - Analysis of memory usage – buffer overflow analysis, system registery
  - Access of executables, mutex
  - Access of disk
  - Replication
  - Attempts to hide other files
  - Attempts to terminate programs
  - Attempts to open other executables.

# DYNAMIC HEURISTIC SCANNING

- Implementation Details
  - Sandbox Environment
  - Reports of malware binaries

- Advantages
  - Can detect encrypted viruses
  - Sandboxing helps in CPU emulation without affecting the system

- Disadvantages
  - CPU Emulation is expensive

# MACHINE LEARNING

# MACHINE LEARNING

- Prototype Extraction – Gonzalez algorithm – O ( kn )
  - K – number of prototypes, n – number of reports
- Clustering using prototypes
  - Complete linkage - $O(k^2 \log k + n)$

  v/s
  - hierarchical clustering - $O(n^2 \log n)$
  - Speed-up factor of square root n/k
- Classification
  - Clustered malware classes used for training and learning behavior
  - O (kn)
- Incremental analysis
  - Better than batch analysis
  - Uses prototypes stored from previous runs
  - $O(nm + k^2 \log k)$
    - M – number of prototypes from previous run

# RESULTS AND ANALYSIS

| | Signature Matching | Static Heuristic Scanning | Dynamic Heuristic Scanning | Artificial Learning |
|---|---|---|---|---|
| Detection of known virus | ✓ | ✓ | ✓ | ✓ |
| Detection of unknown viruses | Fails when signature is unavailable | ✓ | ✓ | **62%** |
| Robustness | Fails when signature is unavailable | ✓ | ✓ ✓ | ✓ |
| False positives | No false positives | ✓ | ✓ | Reduces false positives |
| High speed detection | ✓ ✓ | ✓ | Requires CPU emulation | Learning algorithms consume time |
| Detect metamorphic/ oligomorphic viruses | Fails since virus encrypts itself | ✓ | ✓ ✓ | Efficient only after detection by heuristic scanning |
| Obfuscation | Fails since virus obfuscates itself | ✓ | ✓ ✓ | After heuristic scanning |

# RESULTS AND ANALYSIS

| Technique | F-measure |
|---|---|
| Clustering using Prototype | 0.950 |
| Classification using Prototype | 0.981 |
| Classification using SVM and XML | 0.807 |

$$\text{Precision} = \frac{tp}{tp + fp}$$

$$\text{Recall} = \frac{tp}{tp + fn}$$

$$F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

- Precision - how well individual clusters agree with malware classes
- Recall - extent to which classes are scattered across clusters
- Inverse relation between Precision and Recall
- F-measure – Combines Precision and Recall
  - 1 => perfect classification
  - 0 => completely incorrect classification

# SCREENSHOTS

# References

- Chamorro, E.; Jianchao Han; Beheshti, M.; , "The Design and Implementation of an Antivirus Software Advising System," *Information Technology: New Generations (ITNG), 2012 Ninth International Conference on* , vol., no., pp.612-617, 16-18 April 2012

- Wei Wang; Pengtao Zhang; Ying Tan; Xingui He; , "A Hierarchical Artificial Immune Model for Virus Detection," *Computational Intelligence and Security, 2009. CIS '09. International Conference on* , vol.1, no., pp.1-5, 11-14 Dec. 2009

- Xiao-bin Wang; Guang-yuan Yang, Yi-chao Li, Dan Liu (2008). "Review on the application of artificial intelligence in antivirus detection system," *Proc of 2008 IEEE Conference on Cybernetics and Intelligent Systems,* 506-509.

- Charles P. Pfleeger, Shari Lawrence Pfleeger. "Program Security," in *Security in Computing*, 3rd ed., Prentice Hall, pp.15-67, Dec. 2002.

- *ClamAV User Manual,*2007 - 2011 Sourcefire, Inc. Authors: Tomasz Kojm

- Symantec Corporation, Understanding Heuristics (1997). "Symantec's Bloodhound Technology," *Symantec White Paper Series*.

# Appendix 1 - Implementation of Signature Based Detection

- Use of ClamAV engine

- Provides an API – libClamAV

- Provides an in-memory database of signatures

- Provides regular updates to two database files
  - main.cvd
  - daily.cvd

| | | |
|---|---|---|
| cl_init | Initialize |
| cl_engine | New Engine |
| cl_load | Load Database Directories |
| cl_engine_compile | Build the engine |
| cl_scandesc | Scan file using options |
| cl_free_engine | Free engine |

# Appendix 2 - DEPENDENCY BASED CLASSIFICATION

**Computer Viruses**

- **Computer Architecture Dependency**
- **CPU Dependency**
- **Operating System Dependency**
- **File format Dependency**
- **Date and Time Dependency**
- **Multipartite viruses**

# Appendix 3 - BEHAVIOR BASED CLASSIFICATION

Computer Viruses

- DOS viruses
- Memory Based
- Process viruses
- Kernel viruses

**Memory Based**
- Memory Resident
- Temporary Memory Resident
- Swapping Viruses

**Kernel viruses**
- Boot record viruses
- Master Boot Record Viruses
- Windows NT viruses

# Appendix 4 - HEURISTIC SCANNING PARAMETERS

- Output using sample data
- Unusual destination
- File types and File System
- Memory Usage – Buffer overflow analysis, system registry
- Access of executables
- Access of disk
- Replication
- Attempts to hide other files
- Source code content matched using wild card characters

# Appendix 5 - IMPLEMENTATION OF HEURISTIC SCANNING



| Heuristic Scanning | |
|---|---|
| **Static Heuristic Scanning** | **Dynamic Heuristic Scanning** |
| Scanning without executing | Monitors calls to operating systems |
| Static code analysis | Requires CPU Emulation |
| Exhaustive search of all code snippets not possible | Robust yet time consuming |