# Graph, Entropy and Grid Computing:

## Automatic Comparison of Malware

*Ismael Briones*

*Aitor Gomez*

*PandaLabs Malware Researcher*

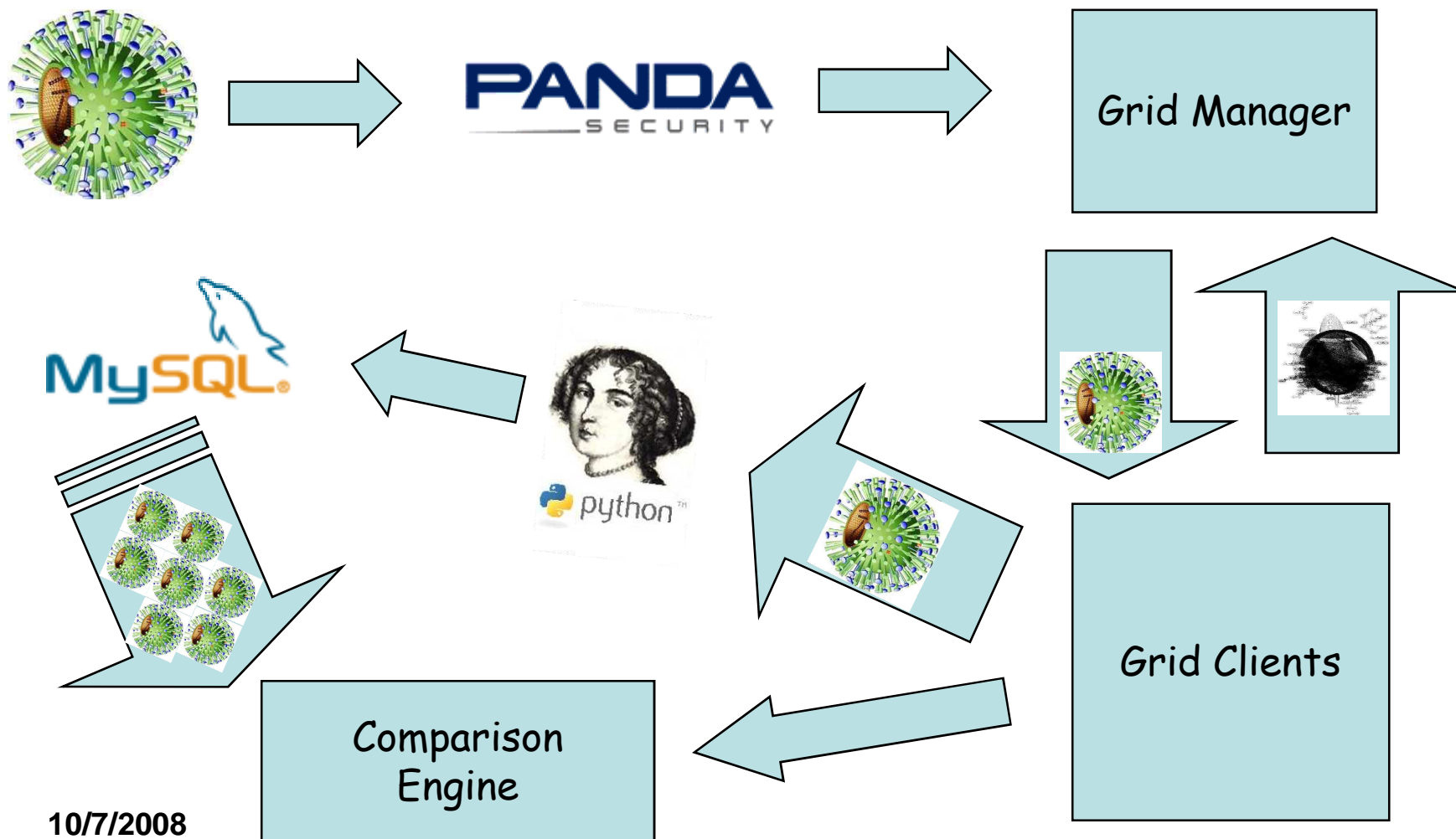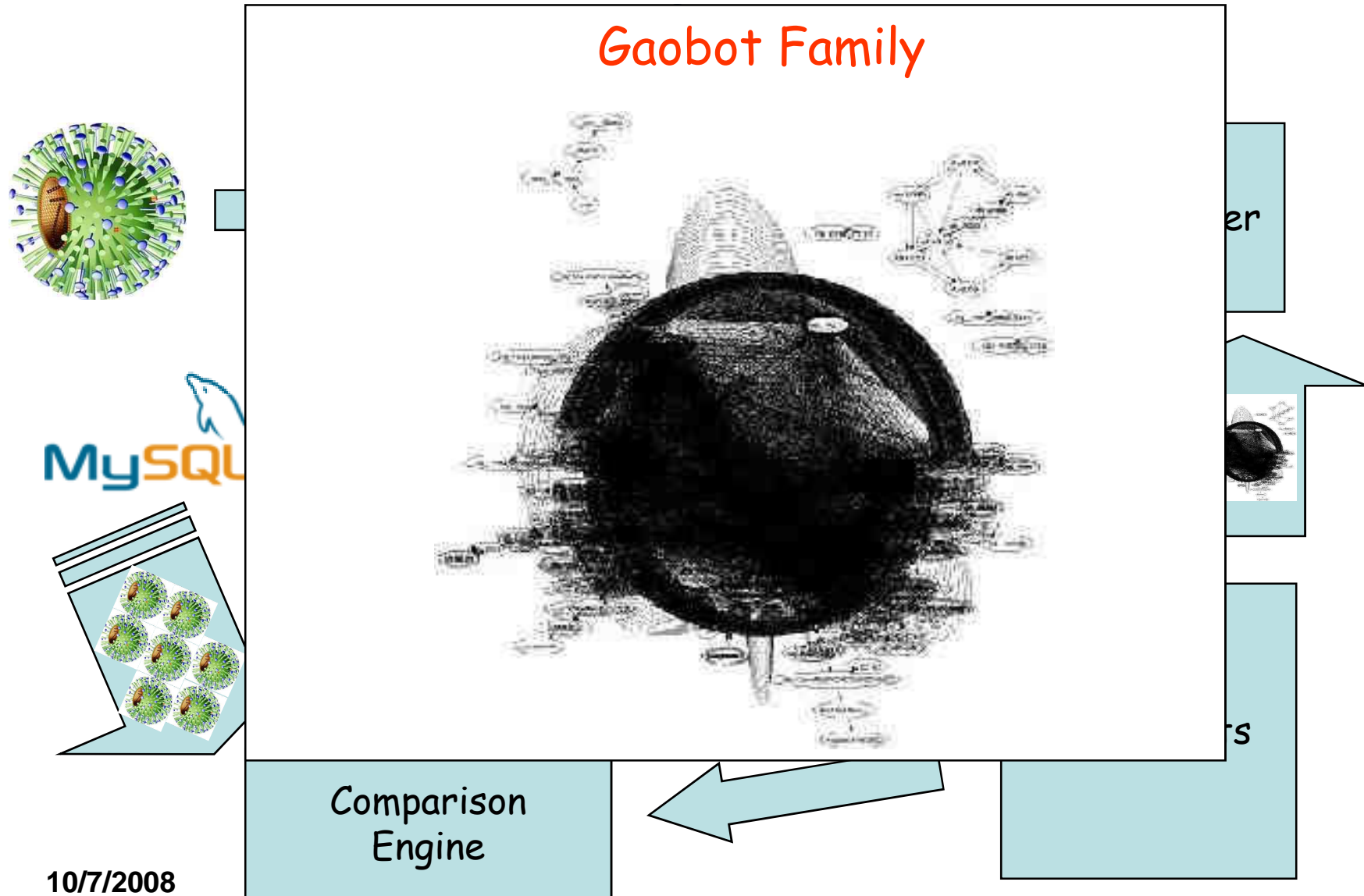Virus Bulletin 2008, Ottawa

**PANDA** SECURITY

# Challenges

➢ Identify new samples

➢ Automatically

➢ ASAP

➢ Improve detection rate

➢ Malware nomenclature

# State of the Art

- ➢ Ero Carrera & Gergely Erdélyi
  - ➢ Digital Genome Mapping
- ➢ Halvar Flake
  - ➢ Lot of researches in graphs analysis
  - ➢ VxClass
- ➢ Marius Gheorghescu
  - ➢ An Automated Virus Classification system

# The System

Grid Manager

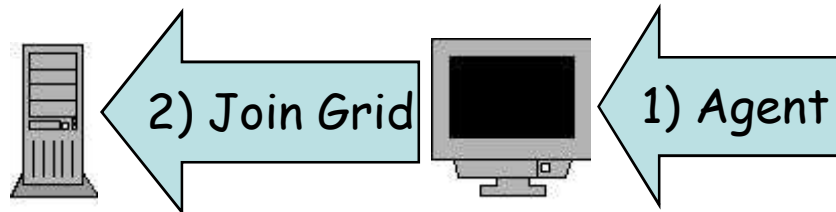Grid Clients

Comparison Engine

Gaobot Family

MySQL

Comparison
Engine

# Grid System

- Analyse as many samples as possible
- IDA Analysis take too much time
- BOINC, GLOBUS, ARC (Advanced Resource Connector)
- XMLRPC

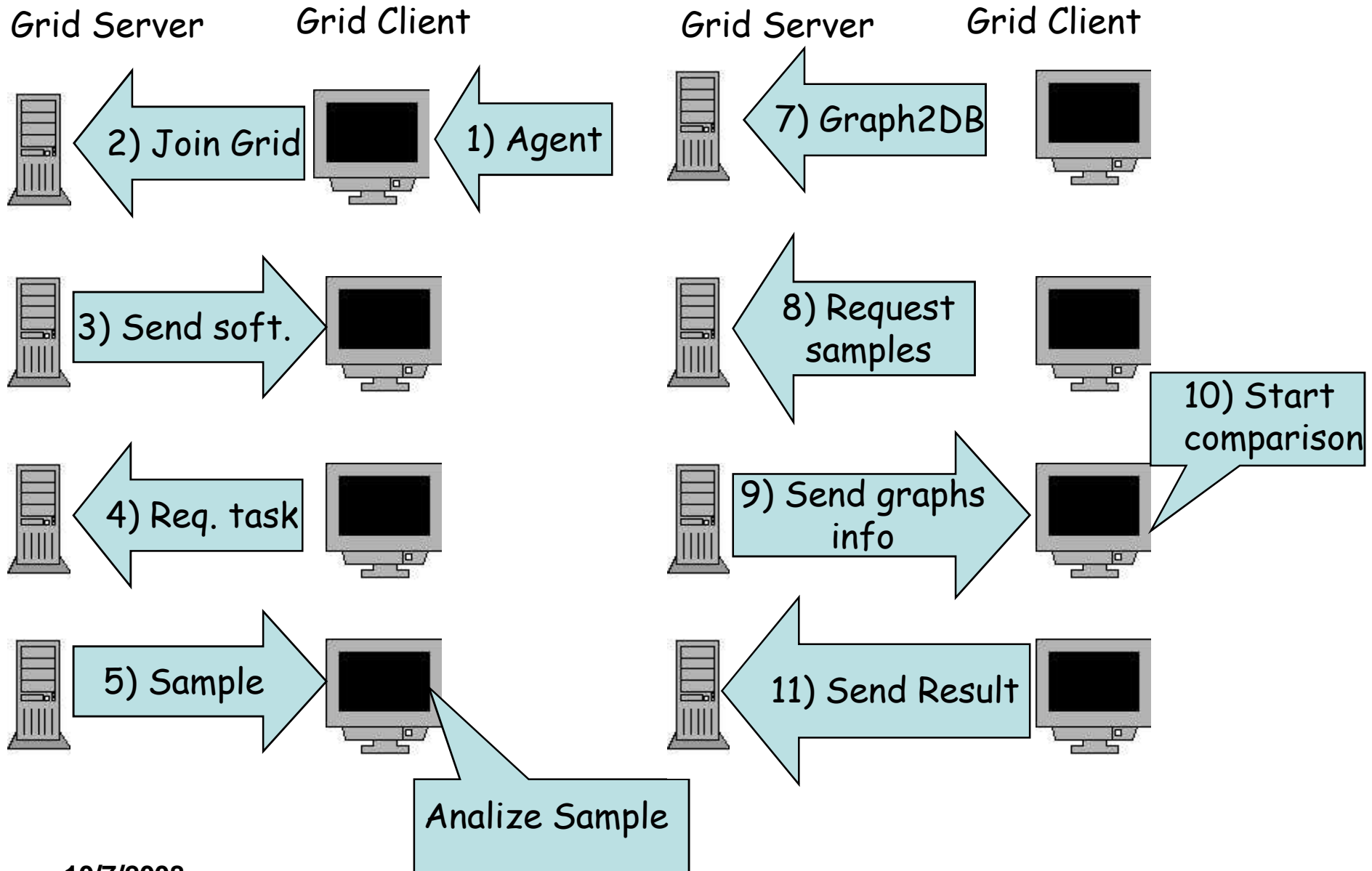# Automatic Comparison of Malware

**Grid Server**   **Grid Client**



2) Join Grid   1) Agent

**Grid System (1.3.3.26)**

Slave Monitor

**Project:** graphStore
**Status:** Initializing....
**Task:** None
**Tasks Done:** 0
**Aborted Tasks:** 0

Registering Slave....
Slave Registered!!

Getting project....

Project Name: graphStore

Project Dir: C:/PandaLabs/Grid/graphStore
Directory "graphStore" exists.

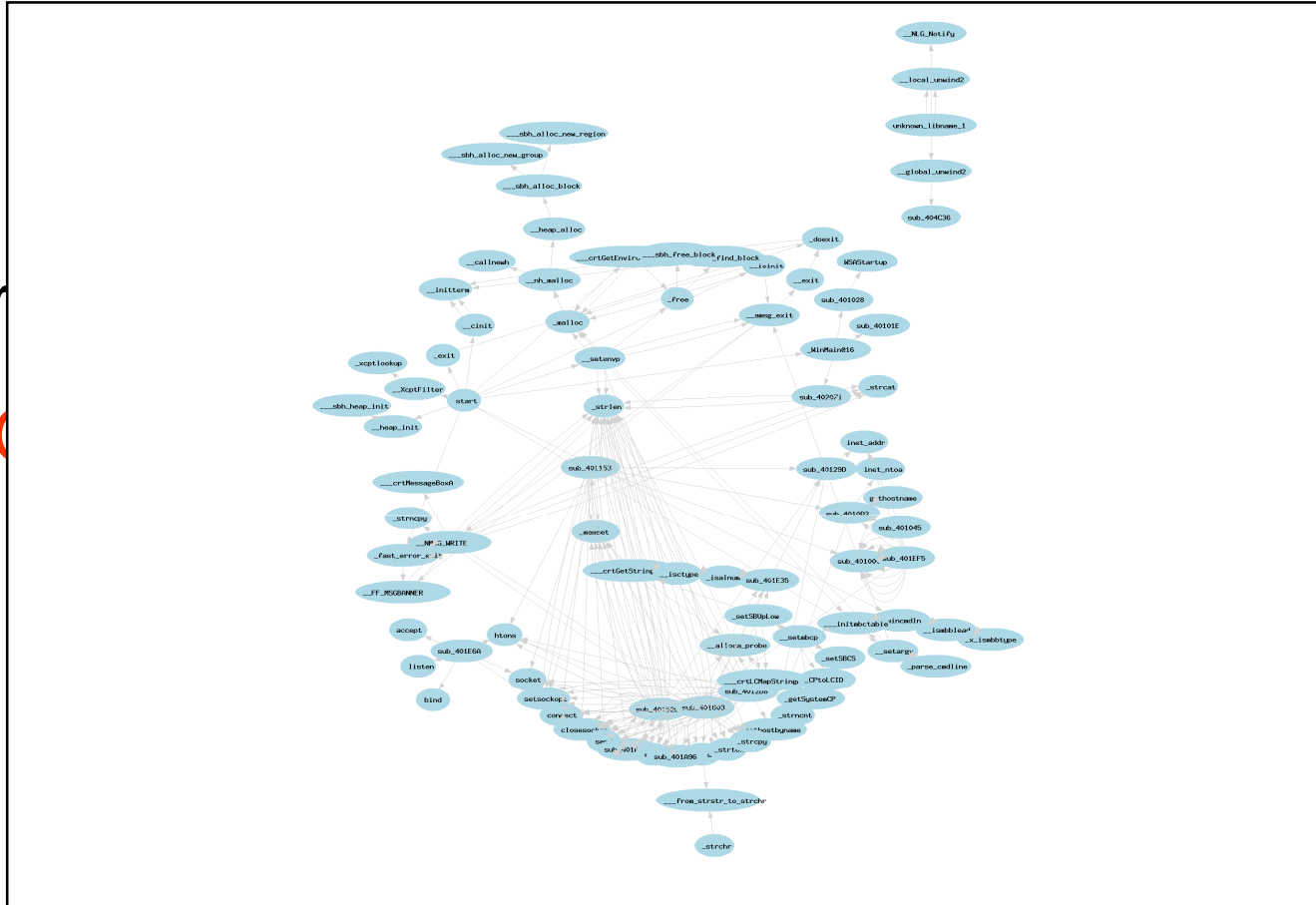Project Software Version: 1.4

# Automatic Comparison of Malware

**PANDA** SECURITY

Grid Server     Grid Client     Grid Server     Grid Client

2) Join Grid

1) Agent

7) Graph2DB

3) Send soft.

8) Request samples

4) Req. task

9) Send graphs info

10) Start comparison

5) Sample

11) Send Result

Analize Sample

**10/7/2008**

**8**

# Sample Analysis

➢Ida Pro + IdaPython

# Sample Analysis

➢Ida Pro + IdaPython

➢Flow Graph
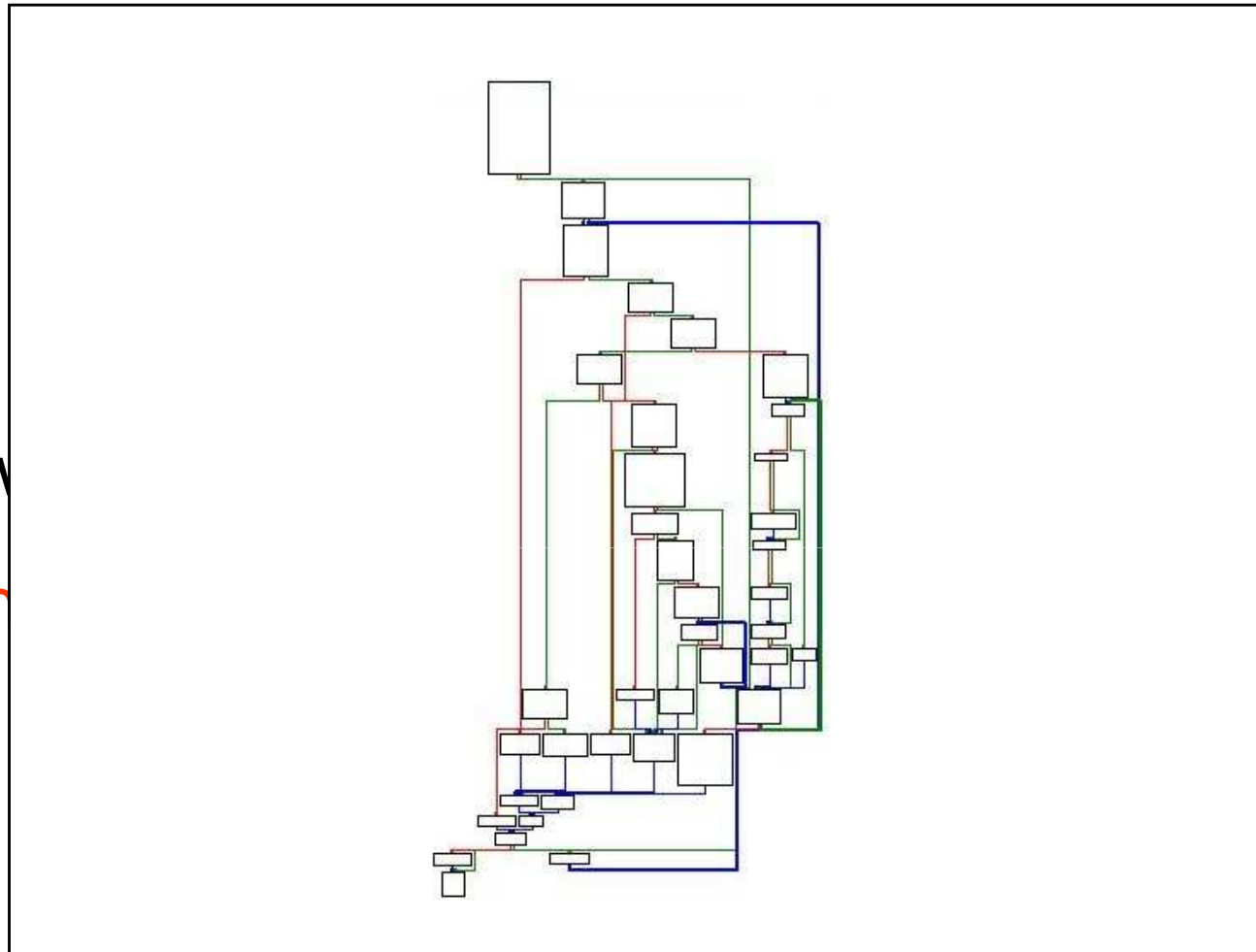
➢Ida Pr

➢Flow

# Adjacency Matrix

➢Ida Pr

➢Flow C

$$\begin{array}{c c c c c c}
 & 1 & 2 & 3 & 4 & 5 \\
1 & 0 & 0 & 1 & 1 & 1 \\
2 & 0 & 0 & 0 & 0 & 1 \\
3 & 0 & 0 & 0 & 0 & 0 \\
4 & 0 & 1 & 0 & 0 & 0 \\
5 & 0 & 0 & 0 & 0 & 0
\end{array}$$

Columns & rows = graph nodes
(i,j) = 1, if i calls j

10/7/2008

12

# Sample Analysis

➢Ida Pro + IdaPython

➢Flow Graph

➢Functions Control Flow Graph (CFG) signature

➢Ida

➢Flow

➢Fun ture

# Automatic Comparison of Malware

# Sample Analysis

➢Ida Pro + IdaPython

➢Flow Graph

➢Functions Control Flow Graph (CFG)

➢Function's crc32

# Sample Analysis

> Ida

> Fl

> Fu

> Fu

```
push    ebp
mov     ebp,
push    ecx
push    ecx
push    esi
call    dword_405044
push    eax
call    sub_40101E
push    1
call    sub_402071
pop     ecx
pop     ecx
call    sub_401028
xor     esi, esi
push    offset aJobaka31 ; "Jobaka31"
push    esi
push    esi
call    dword_405040
```

CRC32(Push+mov+push+push+push+call+…)

# Sample Analysis

➢Ida Pro + IdaPython

➢Flow Graph

➢Functions Control Flow Graph (CFG)

➢Function's crc32

➢ Functions names (sub_, nullsub_,….)

➢Operating Systems and Library Calls (API's)

# Comparison Algorithm

➢**Select best samples**

- Compiler type (Peid Signature)
- File Size
- Number Api Functions
- Number custom functions
- Checksum & Entropy

# Checksum

*'A checksum is a form of redundancy check, […] It works by adding up the basic components of a message, typically the assorted bits* [in our case each byte]*, and storing the resulting value.'*

*[from wikipedia]*
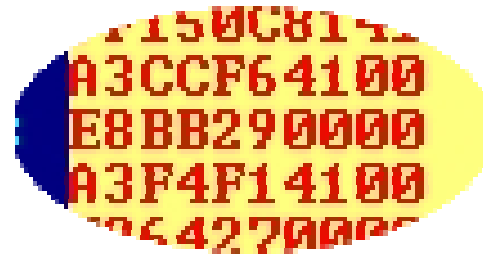
$$(AA)+(BB)+(CC)+(DD) \rightarrow Checksum = 0x30E$$

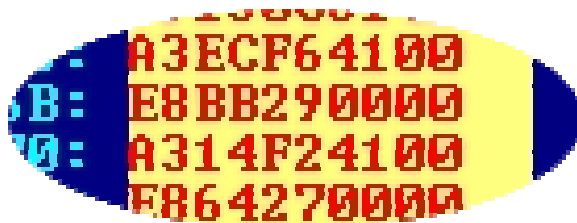# Checksum Properties

➢Similar blocks has very close checksum values

**Checksum: 0x260A**

# Checksum Diff: 0x163

10/7/2008

# Checksum disadvantages

➢A checksum is not changed by:

➢Inserting or deleting zero-valued bytes

$$0xAA+0xBB = 0x165$$
$$0xAA+0x00+0xBB+0x00=0x165$$

➢Reordering of the bytes in the data.

$$0xAA+0xBB = 0x165$$
$$0xBB+0xAA=0x165$$

# Automatic Comparison of Malware

> Checksum value depends on 2 factors:
>> Size of data block
>> Base (basic component: bits, bytes,…)

```
                          255 x 1024 = 0x03FC00
                          (3 bytes)
           256 (1 byte)

1024 bytes     65.536 (2 bytes)   65,535 x 512 =
                                   0x01FFFE00
                                   (4 bytes)

           4.294.976.296 (4 bytes)   4,294,967,295 x
                                     256 =
                                     0xFFFFFFFF00 (5
                                     bytes)
```

➤ Checksum value depends on 2 factors:

➤ Size of data block

➤ Base (basic component: bits, bytes,…)

255 x 1024 = 0x03FC00

Bigger base implies a bigger checksum for the same data block

512 =
00

1024 bytes

4.294.976.296 (4 bytes)

4,294,967,295 x 256 = 0xFFFFFFFF00 (**5 bytes**)

10/7/2008

25

**Bigger base involve a Bigger checksum**

➢Decreases the probability of collision (same checksum)

# Bigger base involve bigger checksum difference

➢ Similar data blocks could get significantly different checksum values

$$Base\ 256\ (BYTE)$$

$$(AA)+(BB)+(CC)+(DD) \rightarrow Checksum = 0x30E$$
$$(BA)+(BB)+(CC)+(DD) \rightarrow Checksum = \underbrace{0x31E}_{(diff=0x10)}$$

$$Base\ 65536\ (WORD)$$

$$(AABB)+(CCDD) \rightarrow Checksum = 0x017798$$
$$(BABB)+(CCDD) \rightarrow Checksum = \underbrace{0x018798}_{(diff=0x1000)}$$

**10/7/2008**

**Bigger base involve bigger checksum difference**

➢ Similar data blocks could get significantly different checksum values

As code is represented as bytes, our approach will use a base of 256 (1byte)

$$(AABB)+(CCDD) \rightarrow Checksum = 0x017798$$
$$(BABB)+(CCDD) \rightarrow Checksum = 0x018798$$
$$(diff = 0x1000)$$

# Our Checksum

➢4KB, in blocks of 1KB, from the beginning of the 1st , 2nd and last sections

| | |
|---|---|
| First Section | 1Kb |
| | 1Kb |
| | 1Kb |
| | 1kb |
| | ... |
| Second Section | 1Kb |
| | 1Kb |
| | 1Kb |
| | 1kb |
| | ... |
| More Sections | ... |
| Last Section | 1Kb |
| | 1Kb |
| | 1Kb |
| | 1kb |
| | ... |

# Checksum Algorithm

➤ Substract checksum from both files

# Automatic Comparison of Malware

Sample A — Sample B

| | | | | |
|---|---|---|---|---|
| **First Section** | 1Kb | **First Section** | 1Kb |
| | 1Kb | | 1Kb |
| | 1Kb | | 1Kb |

ABS(ChksmA – ChksmB) <= EntropyDiff

| | | | |
|---|---|---|---|
| ... | | ... |
| **Second Section** | 1Kb | **Second Section** | 1Kb |
| | 1Kb | | 1Kb |
| | 1Kb | | 1Kb |
| | 1kb | | 1kb |
| | ... | | ... |
| **More Sections** | | **More Sections** | |
| | ... | | ... |
| **Last Section** | 1Kb | **Last Section** | 1Kb |
| | 1Kb | | 1Kb |
| | 1Kb | | 1Kb |
| | 1kb | | 1kb |
| | ... | | ... |

# Automatic Comparison of Malware

**PANDA** SECURITY

Sample A                    Sample B

| First Section | 1Kb | First Section | 1Kb |
|---|---|---|---|
| | 1Kb | | 1Kb |
| | 1Kb | | 1Kb |
| | 1kb | | 1kb |

$$ABS(ChksmA - ChksmB) <= EntropyDiff$$

| Second Section | 1Kb | Second Section | 1Kb |
|---|---|---|---|
| | 1Kb | | 1Kb |
| | 1Kb | | 1Kb |
| | 1kb | | 1kb |
| | ... | | ... |
| More Sections | | More Sections | |
| | ... | | ... |
| Last Section | 1Kb | Last Section | 1Kb |
| | 1Kb | | 1Kb |
| | 1Kb | | 1Kb |
| | 1kb | | 1kb |
| | ... | | ... |

## Sample A

## Sample B

| First Section | 1Kb |
| --- | --- |
| | 1Kb |
| | 1Kb |
| | 1kb |
| | ... |
| Second Section | 1Kb |
| | 1Kb |
| | 1Kb |
| | 1kb |
| | ... |
| More Sect | |
| Last Section | 1Kb |
| | 1Kb |
| | 1Kb |
| | 1kb |
| | ... |

| First Section | 1Kb |
| --- | --- |
| | 1Kb |
| | 1Kb |
| | 1kb |
| | ... |
| Second Section | 1Kb |
| | 1Kb |
| | 1Kb |
| | 1kb |
| | ... |
| Last Section | 1Kb |
| | 1Kb |
| | 1Kb |
| | 1kb |
| | ... |

ABS(ChksmA – ChksmB) <= EntropyDiff

# Checksum vs Entropy

# Checksum Collisions & Entropy

# Checksum Collisions & Entropy



Medium-low entropy
contain code

**10/7/**

37

# Penalized checksum

➢Don't use blocks with high entropy

➢Penalized checksum:

$$if \ (Entropy \geq 7)$$

$$\Downarrow$$

$$PenalizedChecksum = \log2((8 - Entropy)) \ \text{x} \ (-60) \ \text{x} \ Checksum$$

➢Minimal checksum changes in high-entropy blocks will generate distant checksum values.

**10/7/2008**

# Comparison Algorithm

➢**Select best samples**

➢**Start comparison:**

    ➢**Graph Comparison**

**PANDA** SECURITY

➢Matrix A

| | a1 | a2 | a3 | a4 |
|-----|----|----|----|----|
| f1 | 0 | 1 | 0 | 1 |
| f2 | 0 | 0 | 0 | 0 |
| f3 | 1 | 0 | 0 | 0 |
| f4 | 0 | 0 | 1 | 0 |

➢Matrix B

| | a1 | a2 | a3 | a4 |
|------|----|----|----|----|
| f1' | 0 | 1 | 0 | 1 |
| f2' | 0 | 0 | 0 | 0 |
| f3' | 1 | 0 | 0 | 0 |
| f4' | 1 | 0 | 0 | 0 |

**10/7/2008**

**PANDA** SECURITY

➢ **Matrix A**          ➢ **Matrix B**

| | a1 | a2 | a3 | a4 |
|---|---|---|---|---|
| **f1** | 0 | 1 | | 1 |
| **f2** | 0 | 0 | | 0 |
| **f3** | 1 | 0 | | 0 |
| **f4** | 0 | 0 | | 0 |

**Common functions:**
- **API Functions**
- **Same and unique CRC32**
- **Same and unique CFG**

10/7/2008

41

➢Matrix A

| | a1 | a2 | a3 | a4 |
|---|---|---|---|---|
| f1 | 0 | 1 | 0 | 1 |
| f2 | 0 | 0 | 0 | 0 |
| f3 | 1 | 0 | 0 | 0 |
| f4 | 0 | 0 | 1 | 0 |

➢Matrix B

| | a1 | a2 | a3 | a4 |
|---|---|---|---|---|
| f1' | 0 | 1 | 0 | 1 |
| f2' | 0 | 0 | 0 | 0 |
| f3' | 1 | 0 | 0 | 0 |
| f4' | 1 | 0 | 0 | 0 |

## ➢ Matrix A

|    | a1 | a2 | a3 | a4 |
|----|----|----|----|----|
| f1 | 0  | 1  | 0  | 1  |
| f2 | 0  | 0  | 0  | 0  |
| f3 | 1  | 0  | 0  | 0  |
| f4 | 0  | 0  | 1  | 0  |

## ➢ Matrix B

|     | a1 | a2 | a3 | a4 |
|-----|----|----|----|----|
| f1' | 0  | 1  | 0  | 1  |
| f2' | 0  | 0  | 0  | 0  |
| f3' | 1  | 0  | 0  | 0  |
| f4' | 1  | 0  | 0  | 0  |

**f1=f1' and f1 != {f2',f3',f4'}**

## ➤Matrix A

|  | a1 | a2 | a3 | a4 | f1 |
|---|---|---|---|---|---|
| f2 | 0 | 0 | 0 | 0 | 1 |
| f3 | 1 | 0 | 0 | 0 | 1 |
| f4 | 0 | 0 | 1 | 0 | 0 |

## ➤Matrix B

|  | a1 | a2 | a3 | a4 | f1' |
|---|---|---|---|---|---|
| f2' | 0 | 0 | 0 | 0 | 1 |
| f3' | 1 | 0 | 0 | 0 | 1 |
| f4' | 1 | 0 | 0 | 0 | 0 |

**10/7/2008**

➢Matrix A

➢Matrix B

|     | a1 | a2 | a3 | a4 | f1 |
| --- | --- | --- | --- | --- | --- |
| f2  | 0  | 0  | 0  | 0  | 1  |
| f3  | 1  | 0  | 0  | 0  | 1  |
| f4  | 0  | 0  | 1  | 0  | 0  |

|      | a1 | a2 | a3 | a4 | f1' |
| --- | --- | --- | --- | --- | --- |
| f2'  | 0  | 0  | 0  | 0  | 1  |
| f3'  | 1  | 0  | 0  | 0  | 1  |
| f4'  | 1  | 0  | 0  | 0  | 0  |

**f2=f2' and f2 != {f3',f4'}**

# Automatic Comparison of Malware

## ➢ Matrix A

|    | a1 | a2 | a3 | a4 | f1 | f2 |
|----|----|----|----|----|----|----|
| f3 | 1  | 0  | 0  | 0  | 1  | 1  |
| f4 | 0  | 0  | 1  | 0  | 0  | 1  |

## ➢ Matrix B

|     | a1 | a2 | a3 | a4 | f1' | f2' |
|-----|----|----|----|----|-----|-----|
| f3' | 1  | 0  | 0  | 0  | 1   | 0   |
| f4' | 1  | 0  | 0  | 0  | 0   | 1   |

**f3 != {f3',f4'} and f4 != {f3',f4'}**

**Two functions have been identified: f1 and f2**

# Comparison Algorithm

➢ **Select best samples**

➢ **Start comparison:**

    ➢ **Graph Comparison**

➢ **Match more functions with CFG**

# Control Flow Graph

➢CFG signature = 3-tuple vector in Euclidean space

➢Find minimal and unique ED among functions

$$P=(p_x, p_y, p_z)$$
$$Q=(q_x, q_y, q_z)$$

$$\sqrt{(p_x-q_x)^2+(p_y-q_y)^2+(p_z-q_z)^2}$$

**Three-dimensional
Euclidean distance**

# Comparison Algorithm

➢ **Select best samples**

➢ **Start comparison:**

  ➢ **Graph Comparison**

➢ **Match more functions with CFG**

➢ **Index of Similarity**

# Index of Similarity

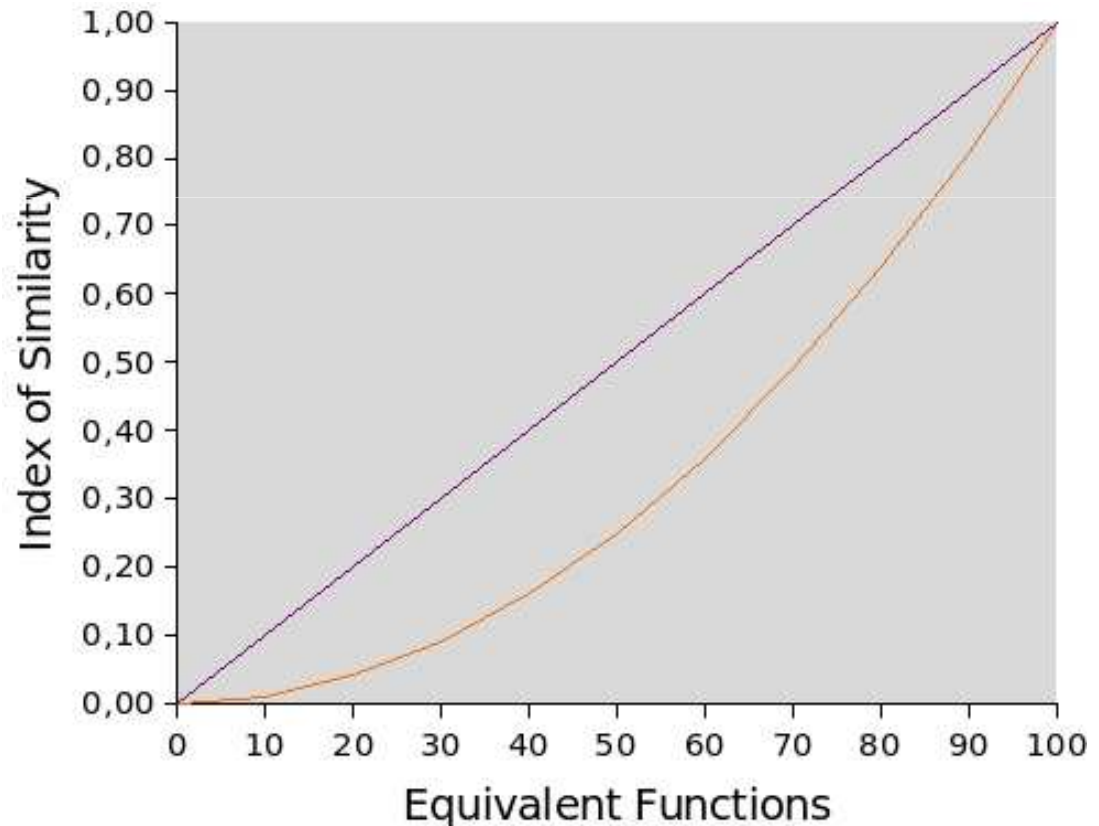➢Measure how close two binaries are.

# Index of Similarity

Index of similarity vs number of matched functions

First Exp.

$$\frac{|A_e||B_e|}{|A_f||B_f|} = \frac{|matchedfunction|^2}{|A_f||B_f|}$$

Second Exp.

$$\frac{|A_e|}{2|A_f|} + \frac{|B_e|}{2|B_f|}$$

# Improvements

➢**More Initial Fixed Points**

# Automatic Comparison of Malware

**PANDA** SECURITY

| More Fixed Points |
| --- |
| SPP (Small Prime Product) |
| Identical String references |
| In/Out degree (similar number of calls to/calls from) |
| Match same name (sub_XXXXXX) if same CRC32 |
| Stack Frame Size (similar stack frame size) |

# Improvements

➢ **More Initial Fixed Points**

➢ **Python version improvements**

# Python version improvements

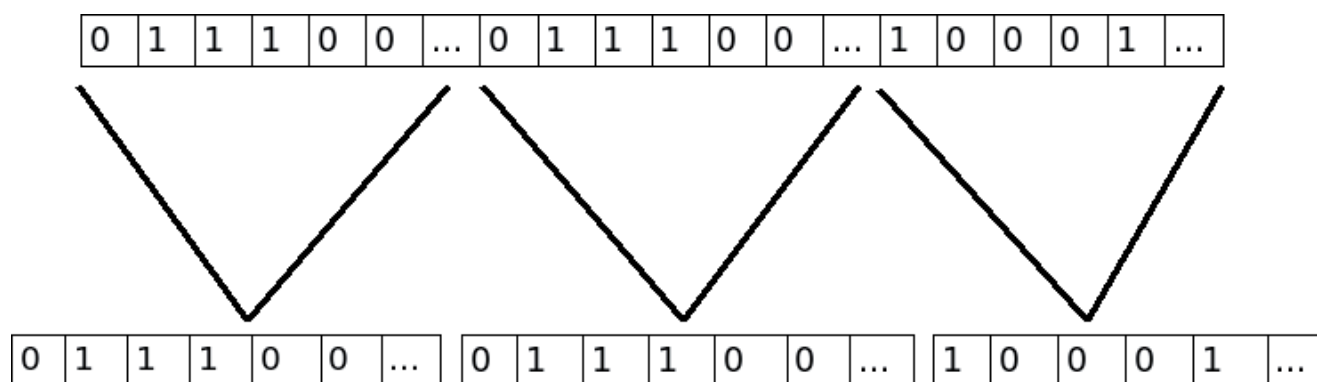➢Psyco, a Just-In-Time compiler

➢Extend Python code with C

➢From Python to C/C++

    ➢Develop the comparison algorithm in C/C++

# Improvements

➢ **More Initial Fixed Points**

➢**Python version improvements**

➢**Matrix rows as bits**

# Matrix rows as bits

➢Avoid string comparison (string comparison is a time-consuming task)

➢Treat rows as groups of bits (100110011101…)

➢Split rows as groups of 32 bits

➢Compare as integers (integer comparison is faster)

# Improvements

➤ **More Initial Fixed Points**

➤ **Python version improvements**

➤ **Matrix rows as bits**

➤ **Stream SIMD extensions**

# Streaming SIMD extensions

➢SSE added eight 128-bits registers: **XMM0-XMM7**

➢Each register packs together four 32-bit integers

➢Compare 4x32 (four integers) in one instruction:

  ➢__m128 _mm_cmpeq_ps(__m128 a, __m128 b)

# Improvements

➢ **More Initial Fixed Points**

➢**Python version improvements**

➢**Matrix rows as bits**

➢**Stream SIMD extensions**

➢**NVIDIA Cuda**

# NVIDIA Cuda

➢Cuda: compiler and SDK for NVIDIA GPUs

➢GPUs:

    ➢Parallel "many-core" architecture

    ➢Each core: thousands of threads simultaneously

➢Not tried yet. Future development:

    ➢Code algorithm for graphics processing unit (GPU)

    ➢Launch one thread for each compared sample

    ➢Launch a thread for each compared row

# Possibilities

➢ **Port information from malware database**

➢**Shared code => generic signature**

    ➢**Functions with same CRC32**

➢**Clusterization of malware automatically**

    ➢**Classify unknown samples (index of similarity)**

# Barriers

➢ **Only unpacked samples**

➢**Hard bound to IDA**

➢**Big database storage**

➢**Delphi and VB samples don't work well**

# Demo

# Thank you very much

# Questions?

ismael.briones@pandasecurity.com
aitor.gomez@pandasecurity.com