



# VxClass

## Clustering Malware, Generating Signatures

zynamics GmbH  
info@zynamics.com



# Overview



- Introduction
- What is VxClass? How does it work?
  - The Malware Pipeline
- Generating signatures from clusters of malware



# Malware

- Keeping up with the flood of malware is hard:
    - Steady increase in number of new variants  
measured by “unique hash” (MD5, SHA1)
    - Off-the-shelf tools to produce malware-variants  
think Swizzor
    - AV-signature databases growing fast  
problems: duplicates, junk, false-positives, ...
    - Time of human malware analysts is scarce  
don't let them do repetitive and error-prone work
- Automated methods needed



# VxClass

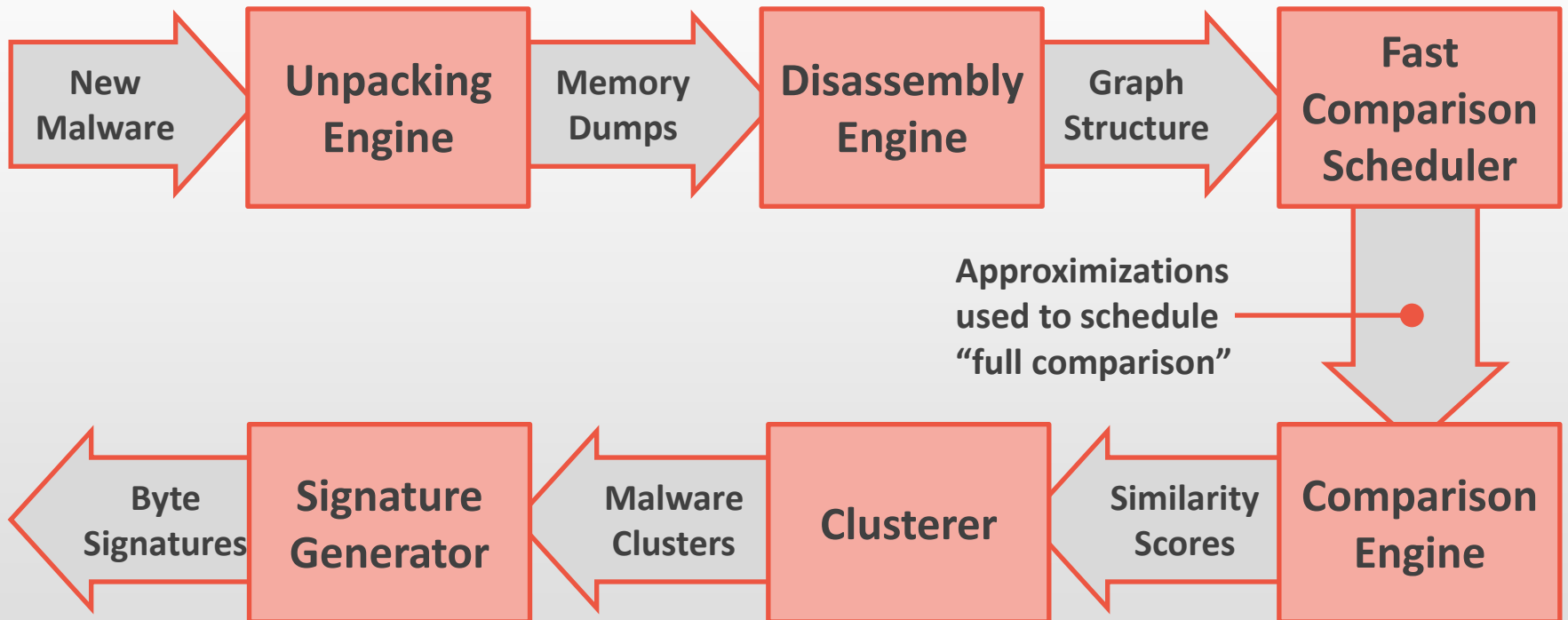


- Full infrastructure for processing new malware
  - Automated generic unpacking to remove crypters based on full-system emulation
  - Comparison engine to detect similarities between executables (based on **zynamics BinDiff**) often better than a human analyst
  - Clustering algorithms for grouping into “families” allows to visualize the malware phylogeny
  - Generation of byte-based AV-signatures for each cluster

# Malware Pipeline



Malware is processed in the following stages





# Unpacking Engine



- Full-system Emulation
  - Uses snapshot of fully booted Win XP SP3 (32-bit)
  - Malware is injected into the system
  - Different execution modes
    - Execute a pre-defined number of instructions
    - Examine repeated snapshots until “good enough”
  - Acquires new processes and new kernel memory  
injected memory into other processes is not yet acquired

Full-system emulation solves many problems

legacy APIs, API detections, ...



# Disassembly Engine



- Processes full executables or memory dumps
- Scan for and recognize typical function prologues
- Generates
  - Flowgraphs
  - Executable Callgraphs
- Compiler and library filtering (FLIRT)

# Comparison Engine



## The “lifeblood” of VxClass

- Based on algorithms developed for the industry-standard **zynamics BinDiff**
- Disregard byte sequences, focus on structural comparison
- Operates on function flowgraphs and the callgraph structure
  - performs comparison based on these structures instead of concrete byte/instruction sequences
- Highly resilient against compiler/platform changes
  - different compilers and settings, even different CPUs (!)





# Comparison Engine

## Example: Mac OS X vs. iPhone

allows cross-CPU comparison

primary

```
8fe19af1  
9af1 lea    eax, (_emergency_mutex -)[ebx]  
9af7 call   __Z20_gthread_mutex_lockP23_opaque_pthread_mutex_t  
9afc mov    edx, ds:(_emergency_used -)[ebx]  
9b02 xor    ecx, ecx  
9b04 cmp    edi, 100h  
9b0a mov    eax, edx  
9b0c jbe    short loc_8FE19B18
```

```
8fe19b18  
9b18 test   al, 1  
9b1a jnz    short loc_8FE19B10
```

```
8fe19b10  
9b10 inc    ecx  
9b11 shr    eax, 1  
9b13 cmp    ecx, 2  
9b16 jz     short loc_8FE19B37
```

```
8fe19b37  
9b37 xor    esi, esi
```

```
8fe19b1c  
9b1c mov    eax, 1  
9b21 shl    eax, cl  
9b23 shl    ecx, 8  
9b26 or    edx, eax  
9b28 lea   esi, (_emergency_buffer  
9b2f mov    ds:(_emergency_used -)  
9b35 jmp    short loc_8FE19B39
```

secondary

```
2fe1746c  
746c LDR    R0, =(_emergency_mutex -)  
7470 MOV    R3, #1  
7474 STR    R3, [SP, #0xA8+var_94]  
7478 ADD    R0, PC, R0  
747c BLX    _pthread_mutex_lock  
7480 LDR    R3, =(_emergency_used -)  
7484 LDR    R1, [PC, R3]  
7488 LDR    R3, [SP, #0xA8+var_9C]  
748c CMP    R3, #0x100  
7490 MOVLS R3, R1  
7494 LDRLS R2, [SP, #0xA8+var_A0]  
7498 BLS    loc_2FE174B0
```

```
2fe174b0  
74b0 TST    R3, #1  
74b4 BNE    loc_2FE174A0
```

```
2fe174a0  
74a0 ADD    R2, R2, #1  
74a4 CHP    R2, #2  
74a8 MOV    R3, R3, LSR#1  
74ac BEQ    loc_2FE174DC
```

```
2fe174dc  
74dc MOV    R3, #0  
74e0 STR    R3, [SP, #0xA8+var_A0]
```

```
2fe174b8  
74b8 MOV    R3, #1  
74bc ORR    R1, R1, R3, LSLR2  
74c0 LDR    R3, =(_emergency_used -)  
74c4 STR    R1, [PC, R3]  
74c8 LDR    R3, =(_emergency_buffer  
74cc ADD    R3, PC, R3  
74d0 ADD    R2, R3, R2, LSL#8  
74d4 STR    R2, [SP, #0xA8+var_A0]  
74d8 B     loc_2FE174E4
```



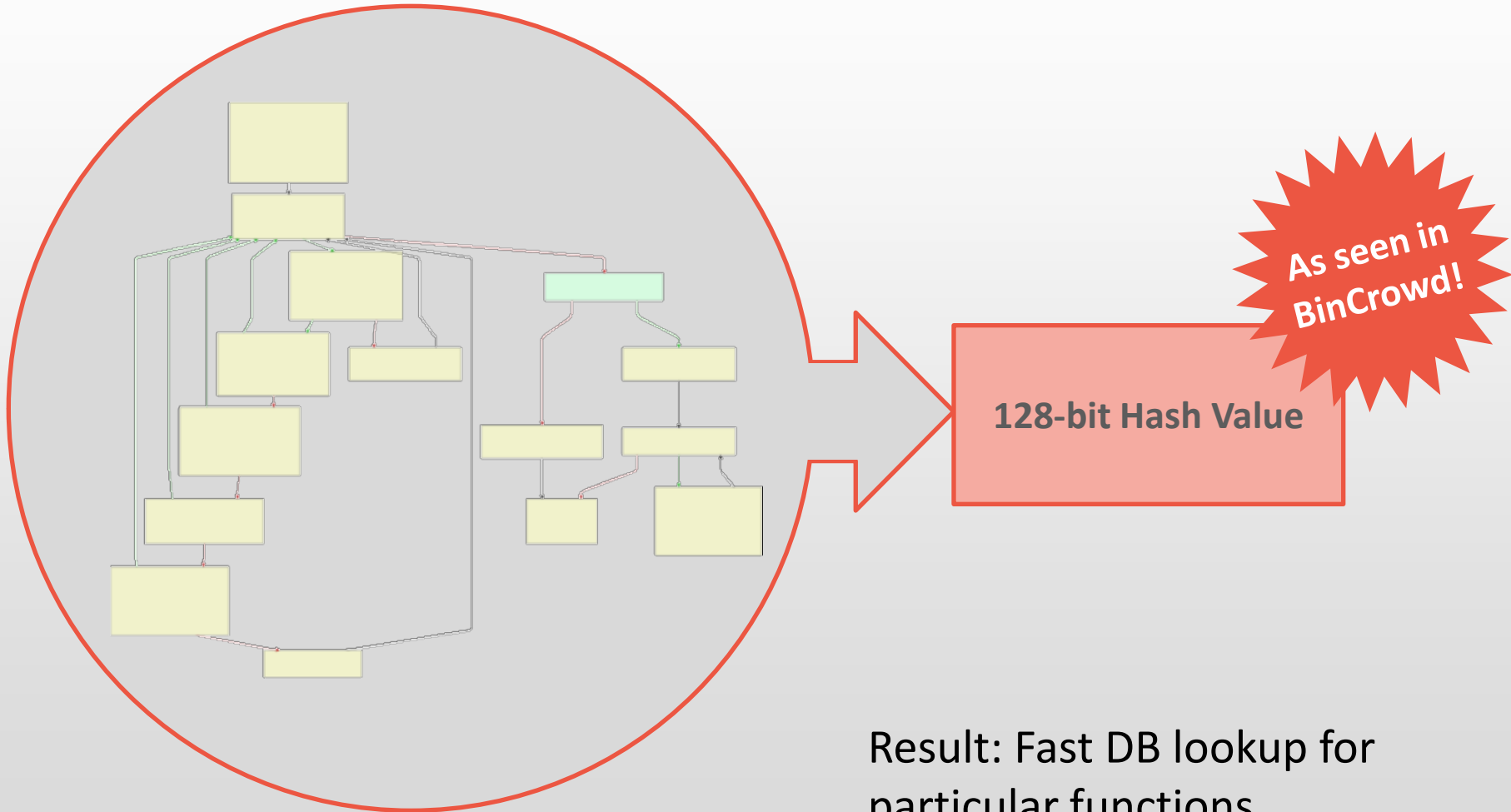
# Scheduling Engine



- Calculating a full similarity matrix is  $O(n^2)$   
prohibitively expensive for large sample sets
- Need to reduce complexity
  - Filter samples that are not similar at all
- Fast comparison engine using  
128-bit flowgraph “hashes” (MD-Index)
  - Calculating MD-Indices is fast
  - Fast comparison via search for common hashes
- Result of fast comparison prioritizes samples  
yields the complete matrix eventually

# MD-Index

Hashing flowgraphs for fast database lookup



Result: Fast DB lookup for particular functions

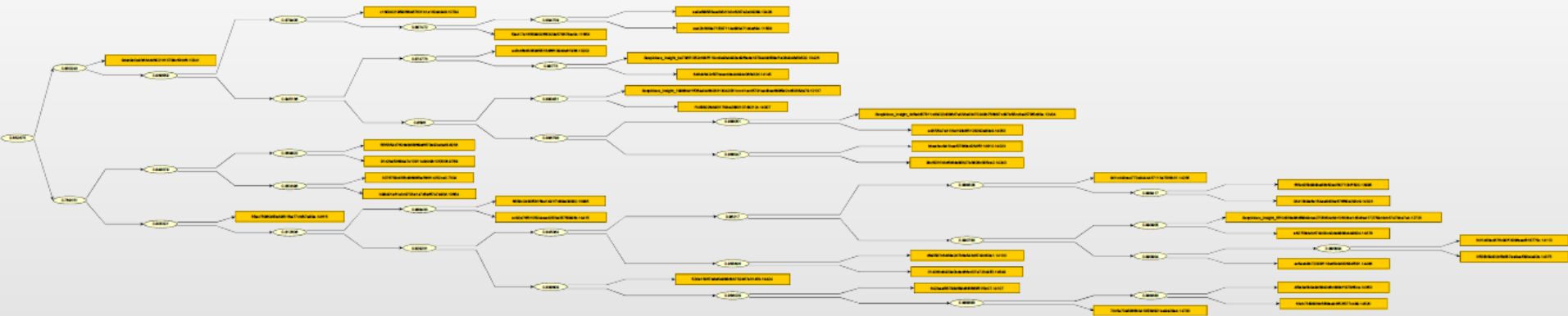


# Clustering

- Based on the similarity matrix generate clusters in different ways
  - Compute connected components
    - use a similarity threshold for graph edges
  - Apply phylogeny algorithms (bioinformatics)
    - yields a family tree
  - Use any other clustering algorithm

# Clustering

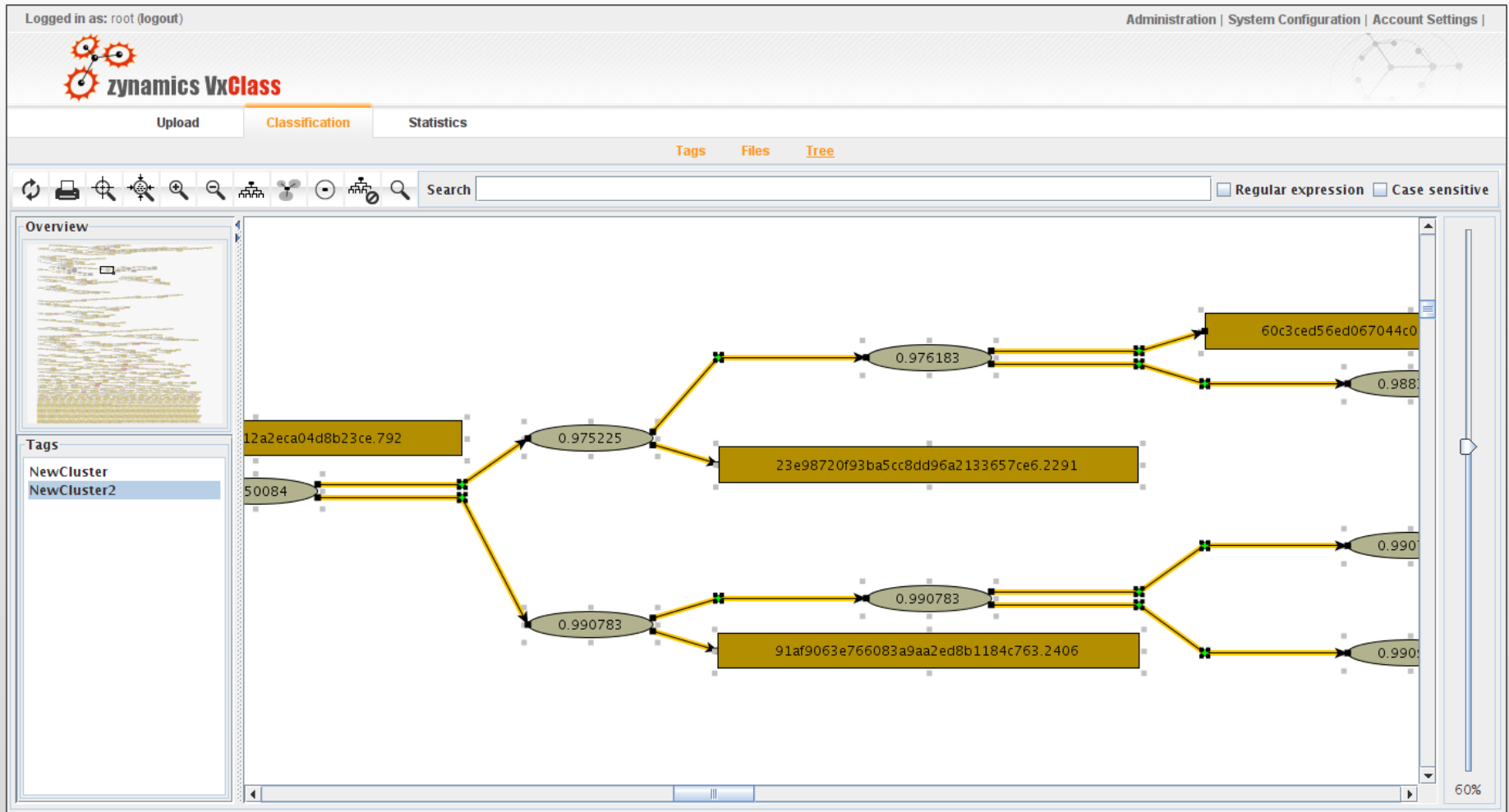
Example tree using phylogeny algorithms



# Clustering

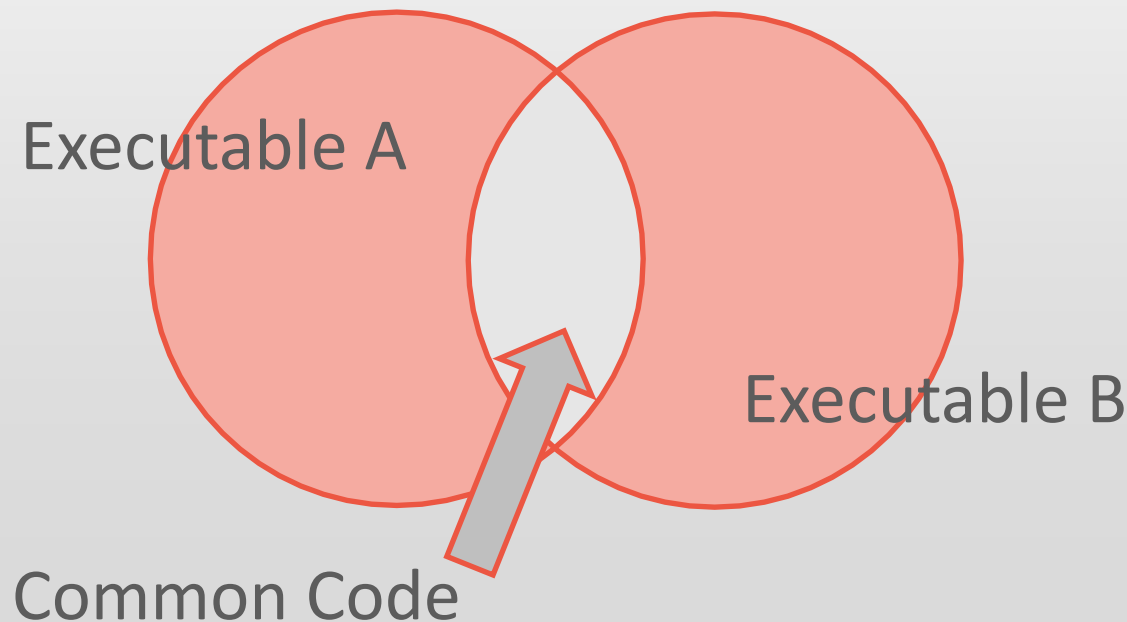
## What to do with those clusters?

tag, name and otherwise analyze them using the web interface



# Signature Generator

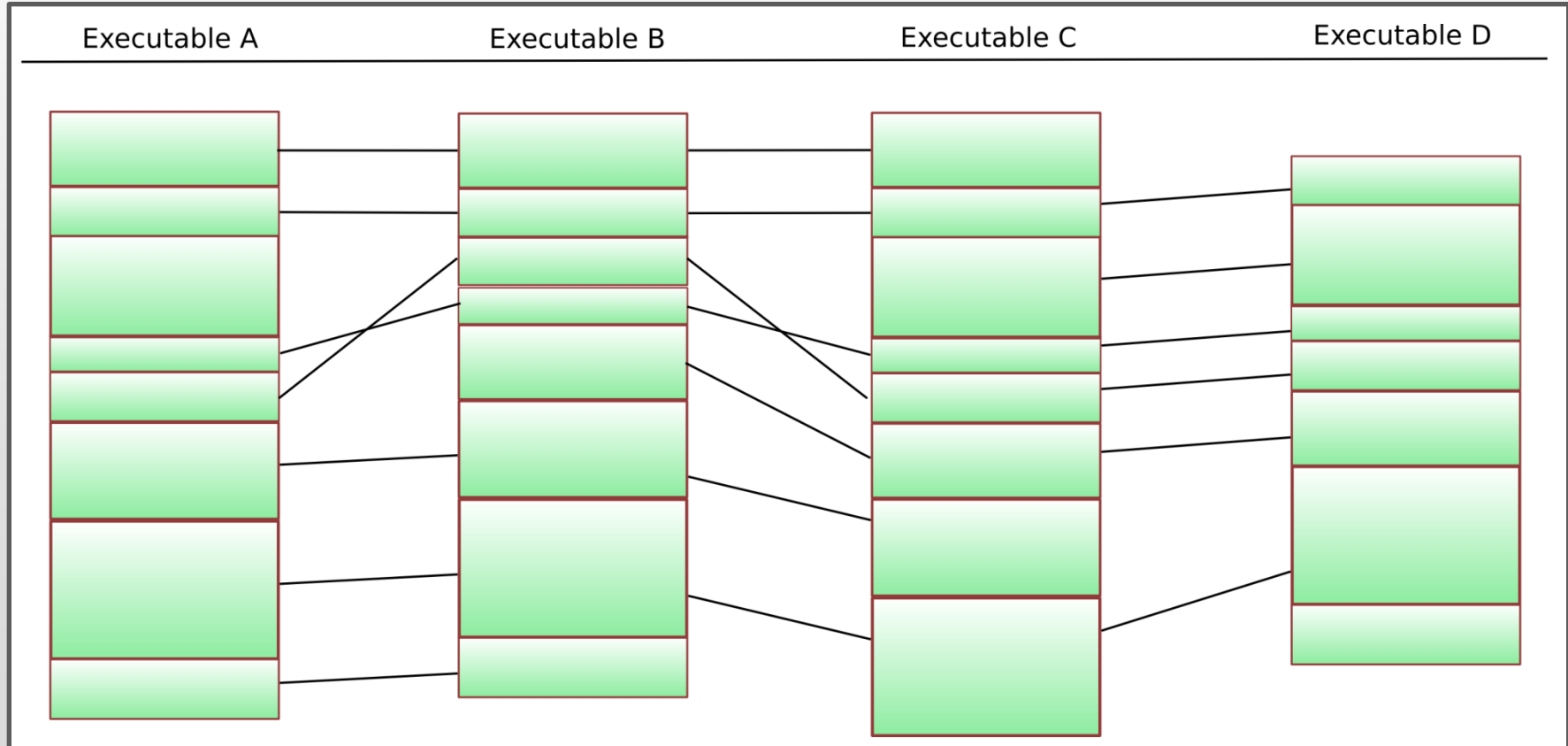
- Fact: Items in the same cluster/family usually share a lot of code
  - Comparison algorithms work like an “intersection operator”:





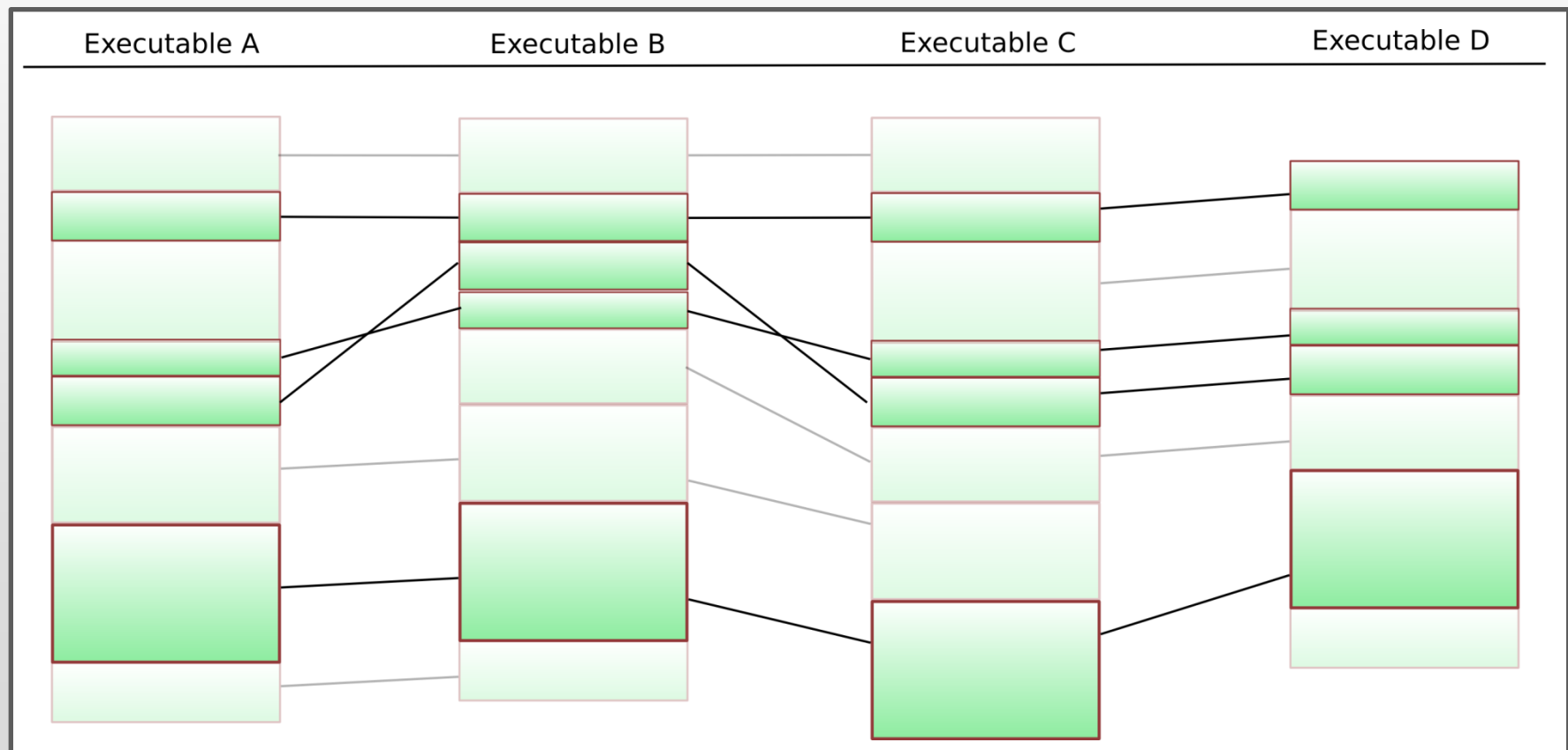
# Signature Generator

- To build a signature, find functions “common to all elements of a cluster”
- Map matched similar functions in each executable:



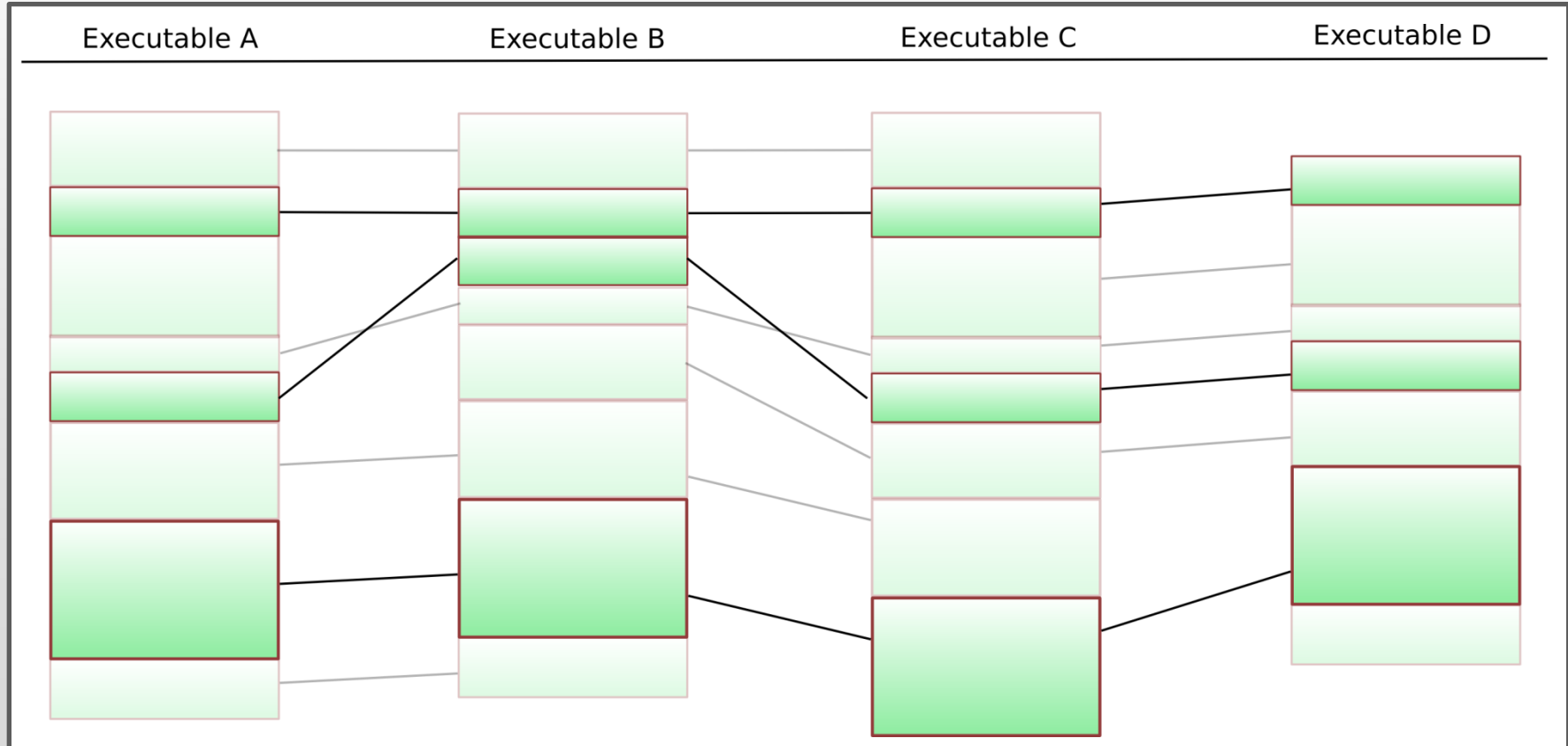
# Signature Generator

- Eliminate functions not present everywhere



# Signature Generator

- A k-LCS algorithm makes sure only functions that are in the same order retained  
this works because functions are identified by their address




# Signature Generator



- Repeat for the basic blocks in each executable
- Results in a “common core” of basic blocks
  - these occur in all items of the cluster and in the same order
- Compute regular k-LCS to determine common byte sequences
  - approximate (exact k-LCS has exponential complexity)
- Fill the “gaps” with wildcards:

```
Worm.SigGen-20080603193253-1876:0:*:03ff*0100*4424*8d54
24*8d4fff*525051*0100*83c418*85c07c*8b4424*8d4c24*4100*
51c64424*02e8*5424*83c40c8d4424*52685c*4100*8b4424*8b46
048b4e08*d1e085c9894604*4100*535657a8018bf175*84c074*8a
d0b9*4100...
```

 Converted to ClamAV format

# Signature Generator



Stats: Classified 5000 random executables from VirusTotal and named resulting clusters.

Signatures were generated and applied to 15000 new executables.

Cluster Name	# executables	sig. size in bytes	new detected variants
Win32.KillAV.Variants	183	1785	111
Win32.Bacuy.Variants	599	27942	863
Win32.SkinTrim	173	290318	356
Win32.SwizzorA	15	69286	929
Win32.WinTrim	114	3925	126
FakeAlert	54	460	0
Win32.Chifrax	12	40098	26



# Signature Generator



- What about false positives?
  - Scanning 22239 known-good executables (ClamAV)
  - Seemingly false positives were always due to bugs
    - bugs in FLIRT, bugs in our library filtering
  - False positive rates empirically around 0.005%
- What about false negatives, then?
  - By construction always either valid signature or none

# Signature Generator



- Signatures consist of malware bytes minus the “variable” bits
  - generated signatures carry some “predictive power”
  - All except one of the generated signatures caught some “new” variant of the same malware



# Performance

- One (beefy) machine processes  $\approx 1400$  samples/day
  - Includes unpacking
  - Higher performance if specific unpacking happens first
- Scalable: VxClass routinely runs on a 4-machine cluster
  - Scaling to 20-25 machines should be possible





# Limitations



- Heavy obfuscation of control flow
  - breaks classification and leads to empty signatures
- Virtualizing packers
- Unpacking only works on 32-bit Windows
  - No Linux/Mac OS X/Mobile unpacking
  - 64-bit support is in the works
- Using the generated signatures in a AV product requires good unpacking capability
  - signatures are generated post-unpacking

But: Manual intervention possible (upload IDBs)



# Questions?



<http://zynamics.com/>

<http://blog.zynamics.com/>

[christian.blichmann@zynamics.com](mailto:christian.blichmann@zynamics.com)