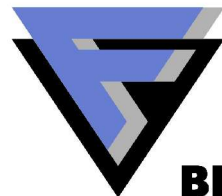


# Digital Genome Mapping Advanced Binary Analysis

Ero Carrera

Gergely Erdélyi

**F-SECURE<sup>®</sup>**



**BE SURE.**

# THE MAP

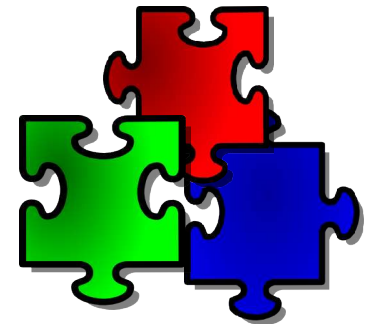
- Introduction
- IDA Pro + Python
- IDAPython Demo
- Introducing graphs
- Creating graphs
- Using graphs

● YOU ARE HERE



# Challenges in Security Research

- Overwhelming number of samples to check
- Size of code to be analysed
- Repetitive analysis of similar samples
- Difficulties of classification
- Use of non-intuitive tools



# IDA Pro

- Full name is Interactive DisAssembler Pro
- Advanced disassembler
- Supports more than 30 processors
- Supports dozens of executable file types
- Has library recognition (FLIRT)
- Extensible by scripts and plugins



# Extending IDA Pro

- Using IDC
  - C-like script language
  - Dynamically typed
  - Lacks higher level construct and data types
- IDA plugins
  - Written in C/C++
  - Compiled and linked into binary code
  - Loaded as DLLs



```
#include <idc.idc>
```

```
static main ()
```

```
{
```

```
    auto baseaddr;
```

```
    auto key;
```

```
    auto c;
```

```
    baseaddr = ScreenEA();
```

```
    key = 0xff;
```

```
    while (1) {
```

```
        c = Byte(baseaddr);
```

```
        if (c == 0) {
```

```
            break;
```

```
        }
```

```
        PatchByte(baseaddr, c ^ key);
```

```
        baseaddr++;
```

```
    }
```

```
}
```



# Extension Language Requirements

- Easy to learn, remember and read
- Supports easy refactoring
- User error tolerant
- Comes with a complete environment
- Provides complex data structures
- Open source
- Free (as in Free beer and Freedom)



# Python

- Interpreted
  - Interactive
  - Object-oriented
  - Dynamically typed
  - Supports modules
  - Has classes and exceptions
  - Very high level data types
    - lists
    - tuples
    - hashes (dictionaries)
  - Extensive libraries
    - OS support
    - Threading
    - Data processing
    - Networking
  - Runs on...
    - Linux, \*NIX
    - Windows, OS/2
    - Amiga, etc
- ...and a lot more

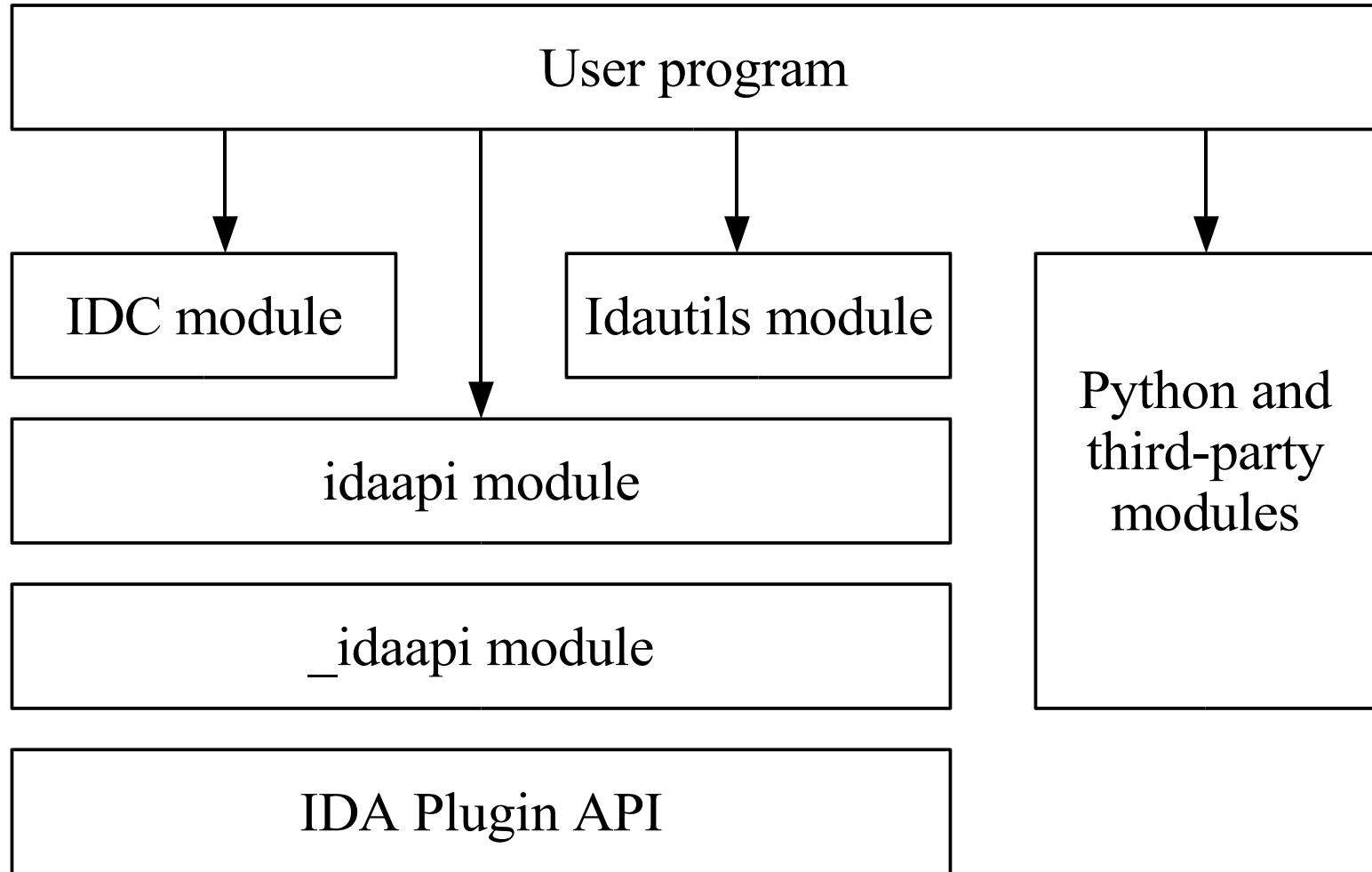


# IDA + Python = IDAPython

- IDAPython is a plugin for IDA Pro
- Provides access for most of IDA Plugin API
  - Functions wrapped with SWIG
  - IDA Plugin API is provided in the module '\_idaapi'
- Python libraries and extensions can be utilized
- IDAPython has an IDC compatibility module
- The 'idautils' module provides abstract utility functions



# Structure of IDAPython



```
#include <idc.idc>
```

```
static main()
```

```
{
```

```
    auto ea, func, ref;
```

```
    // Get current ea
```

```
    ea = ScreenEA();
```

```
    // Loop from start to end in the current segment
```

```
    for (func=SegStart(ea);
```

```
        func != BADADDR && func < SegEnd(ea);
```

```
        func=NextFunction(func))
```

```
    {
```

```
        // If the current address is function process it
```

```
        if (GetFunctionFlags(func) != -1)
```

```
        {
```

```
            Message("Function %s at 0x%x\n", GetFunctionName(func), func);
```

```
            // Find all code references to func
```

```
            for (ref=RfirstB(func); ref != BADADDR; ref=RnextB(func, ref))
```

```
            {
```

```
                Message("    called from %s(0x%x)\n", GetFunctionName(ref), ref);
```

```
            }
```

```
        }
```

```
    }
```

```
}
```

PLAIN IDC



```
from idaapi import *

# Get current ea
ea = get_screen_ea()

# Get segment class
seg = getseg(ea)

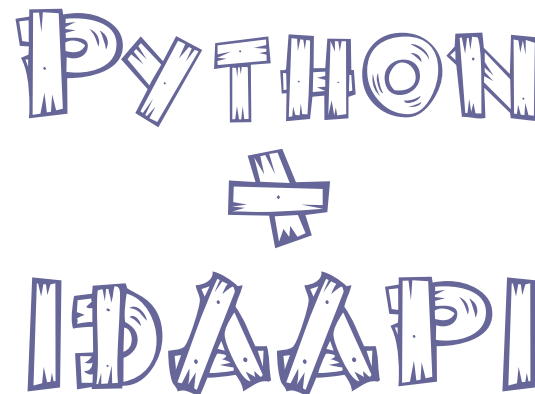
# Loop from segment start to end
func = get_func(seg.startEA)

while func != None and func.startEA < seg.endEA:
    funcea = func.startEA
    print "Function %s at 0x%x" % (GetFunctionName(funcea), funcea)

    ref = get_first_cref_to(funcea)

    while ref != BADADDR:
        print "  called from %s(0x%x)" % (get_func_name(ref), ref)
        ref = get_next_cref_to(funcea, ref)

    func = get_next_func(funcea)
```



PYTHON



IDAUTILS

```
from idautils import *
```

```
# Get current ea  
ea = ScreenEA()
```

```
# Loop from start to end in the current segment  
for funcea in Functions(SegStart(ea), SegEnd(ea)):  
    print "Function %s at 0x%x" %  
        (GetFunctionName(funcea), funcea)
```

```
# Find all code references to funcea  
for ref in CodeRefsTo(funcea, 1):  
    print "    called from %s(0x%x)" %  
        (GetFunctionName(ref), ref)
```



# THE MAP

- Introduction
- IDA Pro + Python
- IDAPython Demo
- Introducing graphs
- Creating graphs
- Using graphs

● YOU ARE HERE



# Availability

IDAPython distribution site:

<http://www.d-dome.net/idapython/>



# THE MAP

- Introduction
- IDA Pro + Python
- IDAPython Demo
- Introducing graphs
- Creating graphs
- Using graphs

● YOU ARE HERE



# Introducing Graphs

- Visual output
  - Graphs provide ways of visualizing complex structures.
  - Human brain is good at capturing patterns, computer has to be taught
- Analysis playground
  - Algorithms can be applied to graphs









# THE MAP

- Introduction
- IDA Pro + Python
- IDAPython Demo
- Introducing graphs
- **Creating graphs**
- Using graphs

● **YOU ARE HERE**



# Creating Graphs

- What can be show in a graph?
  - Just *Nodes* and *Edges*
- How do we represent a virus in a graph?
  - Viruses are like any other program
  - Programs are generally composed of functions
  - Functions are individual parts, used as building blocks
  - We could show the relationships between them



# Infrastructure

In order to represent them we created the following software infrastructure:

- Data exporting (From IDA DB into an XML format)
- Data loader (Parses the data into easily manipulable objects)
- Methods of outputting such objects and their relations as graphs



# Exporting data

An IDC script exports information from IDA into XML:

- Every defined function and instruction
- The full binary representation
- Mnemonic and operands as shown by IDA
- All the defined referenced strings and names
- All code references
- Code flags: whether the code is flow
- Miscellaneous flags: whether a function belongs to a known library or is a thunk



```
<function name="" is_library="true" is_thunk="true">
<i at="401000" is_code="true"><d>31c0</d><mnem>xor</mnem><o
n="0" t="1">eax</o><o n="1" t="1">eax</o><ref_data
from="40102a"/></i>
[...]
<i at="40100e" is_code="true"
is_flow="true"><d>740f</d><mnem>jz</mnem><o n="0"
t="7">4198431</o></i>
[...]
<i at="40101a" is_code="true"
is_flow="true"><d>b803000000</d><mnem>mov</mnem><o n="0"
t="1">eax</o><o n="1" t="5">3</o></i>
<i at="40101f" is_code="true"
is_flow="true"><d>c3</d><mnem>retn</mnem><ref
from="40100e"/></i>
```

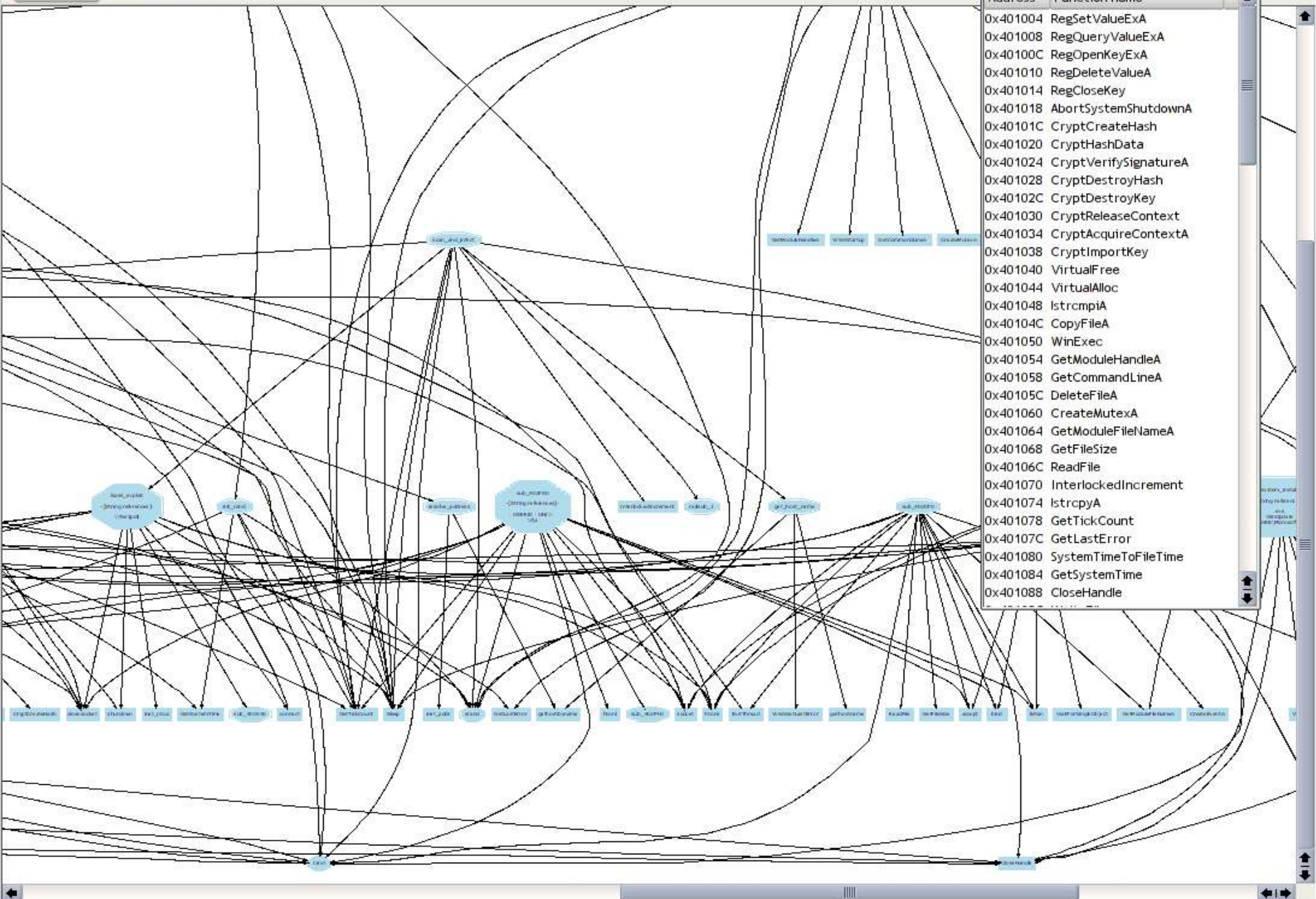


# Data loader

- REML Python loader
  - Provides an interface to the data in a binary, such as functions; allowing to handle them as objects
- Graph handling tools

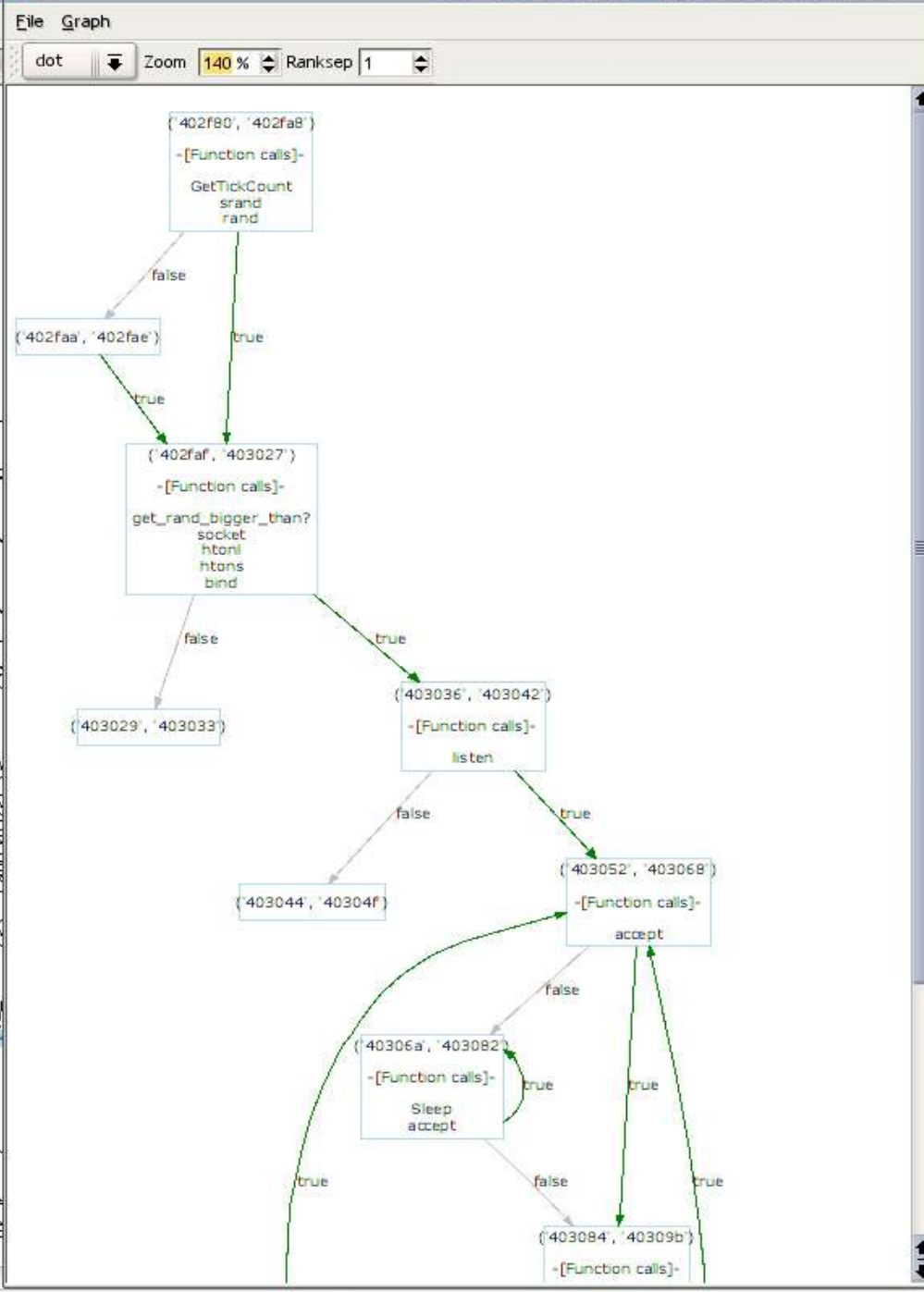


dot Zoom 100 % Ranksep 1



File Graph Windows

dot Zoom 100 %

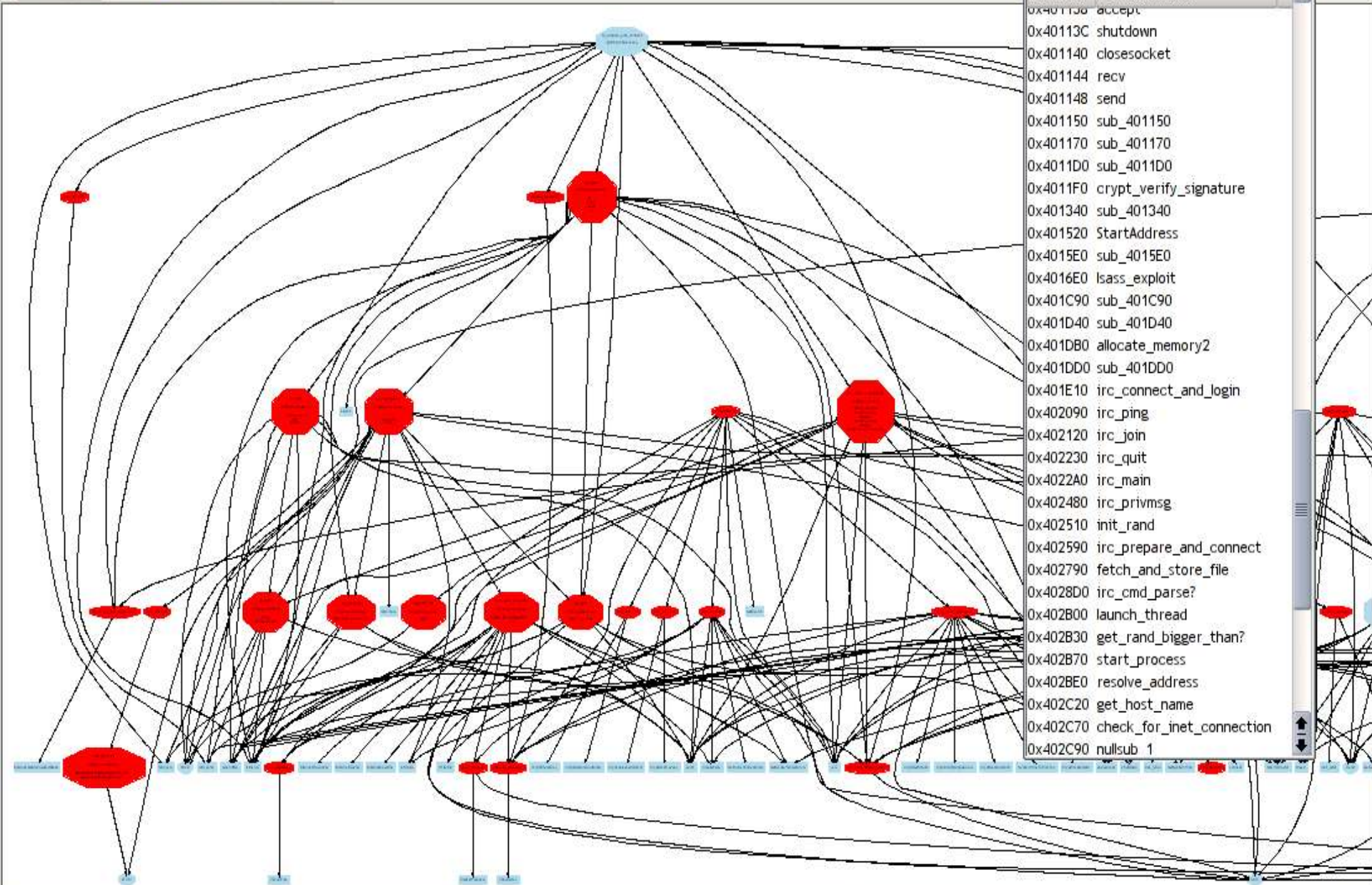


reveng\_browser.py

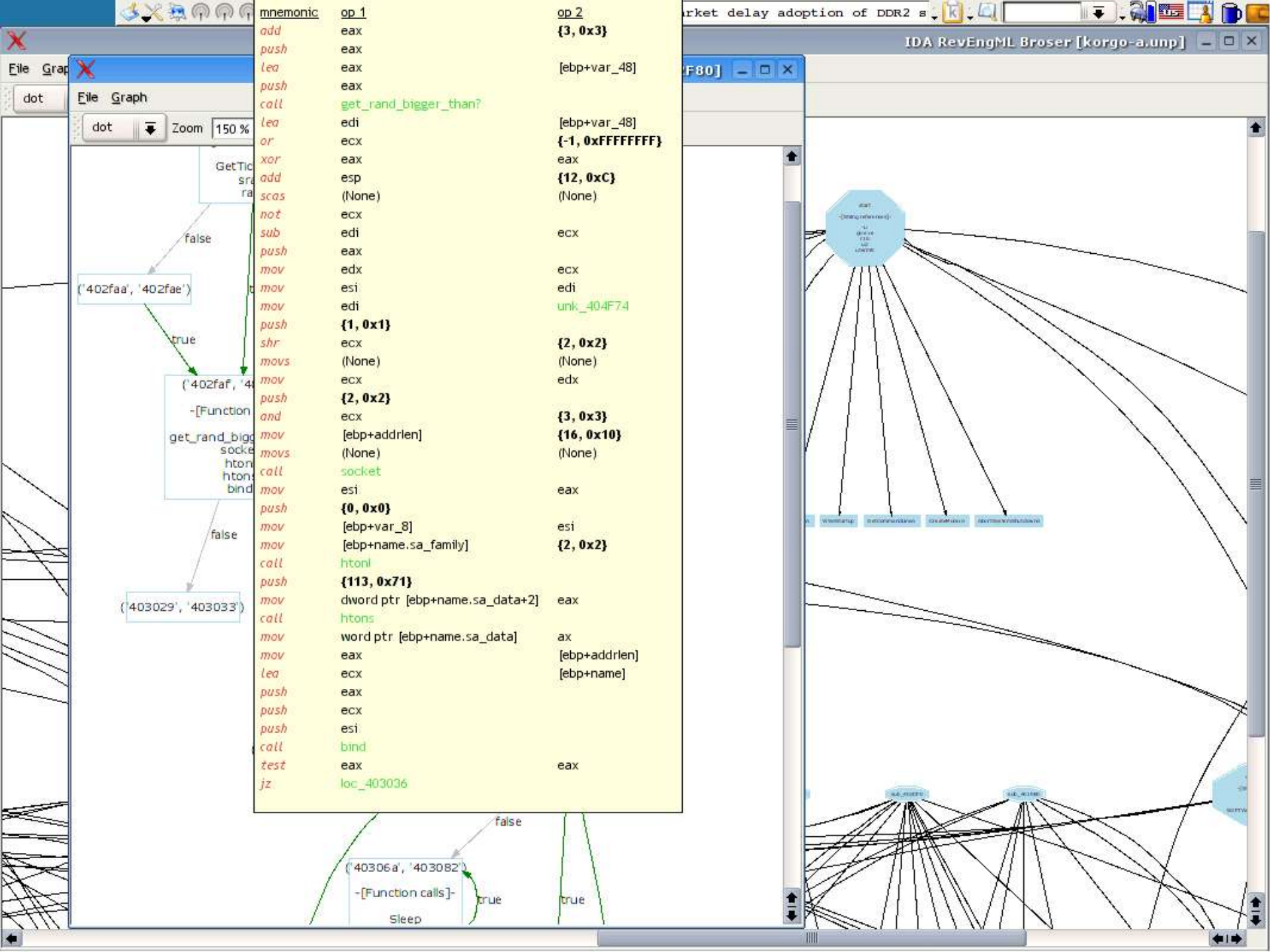
Address	Function name
0x401004	RegSetValueExA
0x401008	RegQueryValueExA
0x40100C	RegOpenKeyExA
0x401010	RegDeleteValueA
0x401014	RegCloseKey
0x401018	AbortSystemShutdownA
0x40101C	CryptCreateHash
0x401020	CryptHashData
0x401024	CryptVerifySignatureA
0x401028	CryptDestroyHash
0x40102C	CryptDestroyKey
0x401030	CryptReleaseContext
0x401034	CryptAcquireContextA
0x401038	CryptImportKey
0x401040	VirtualFree
0x401044	VirtualAlloc
0x401048	IstrcmpA
0x40104C	CopyFileA
0x401050	WinExec
0x401054	GetModuleHandleA
0x401058	GetCommandLineA
0x40105C	DeleteFileA
0x401060	CreateMutexA
0x401064	GetModuleFileNameA
0x401068	GetFileSize
0x40106C	ReadFile
0x401070	InterlockedIncrement
0x401074	IstrcpyA
0x401078	GetTickCount
0x40107C	GetLastError
0x401080	SystemTimeToFileTime
0x401084	GetSystemTime
0x401088	CloseHandle

File Graph Windows

dot Zoom 60% Ranksep 1



Address	Function name
0x401130	accept
0x40113C	shutdown
0x401140	closesocket
0x401144	recv
0x401148	send
0x401150	sub_401150
0x401170	sub_401170
0x4011D0	sub_4011D0
0x4011F0	crypt_verify_signature
0x401340	sub_401340
0x401520	StartAddress
0x4015E0	sub_4015E0
0x4016E0	isass_exploit
0x401C90	sub_401C90
0x401D40	sub_401D40
0x401DB0	allocate_memory2
0x401DD0	sub_401DD0
0x401E10	irc_connect_and_login
0x402090	irc_ping
0x402120	irc_join
0x402230	irc_quit
0x4022A0	irc_main
0x402480	irc_privmsg
0x402510	init_rand
0x402590	irc_prepare_and_connect
0x402790	fetch_and_store_file
0x4028D0	irc_cmd_parse?
0x402B00	launch_thread
0x402B30	get_rand_bigger_than?
0x402B70	start_process
0x402BE0	resolve_address
0x402C20	get_host_name
0x402C70	check_for_inet_connection
0x402C90	nullsub_1



# THE MAP

- Introduction
- IDA Pro + Python
- IDAPython Demo
- Introducing graphs
- Creating graphs
- Using graphs

● YOU ARE HERE



# Comparison

Our work usually faces the need to identify and classify malware into families.

Members of a family, must share some similarities. Those become quite evident in graphs of their call trees.



# Binary comparison, the algorithm

We want to map functions from one binary to another:

- We use the Windows API as a ground set for matching
- We match the API callers and iteratively ascend the graph



# Signatures

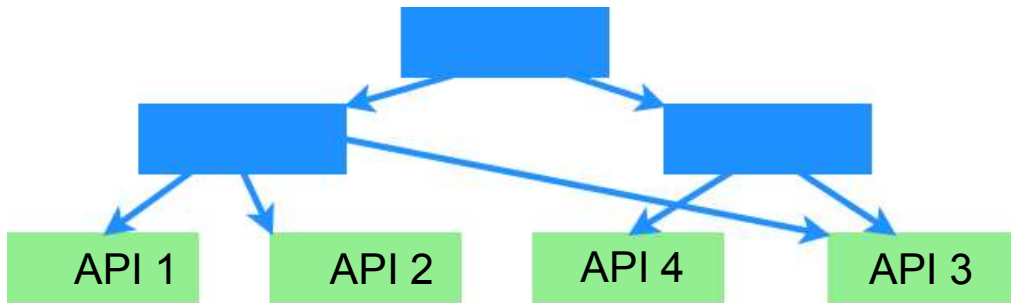
For deciding whether two functions match, a signature must to agree, we use two approaches

- Call Tree Signatures. Giving an idea of the connectivity of a function within the malware body
- CFG signatures. Relying only on the internal structure of a function

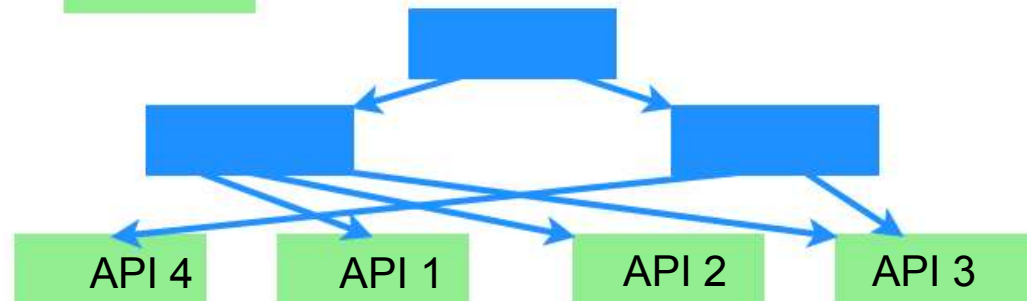


# Call Tree Signatures

Call Tree Signatures



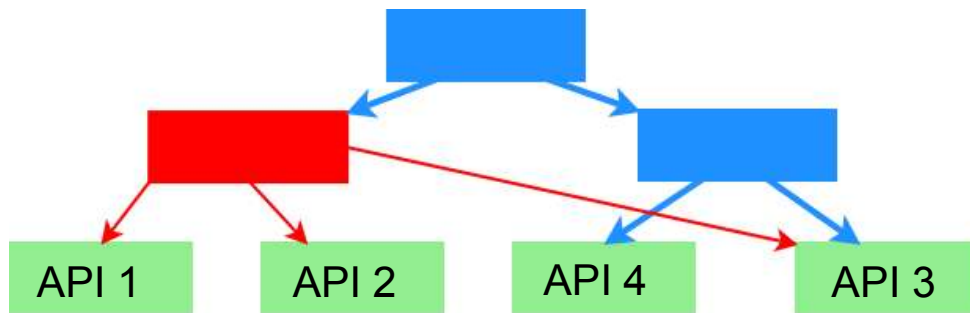
Sample A



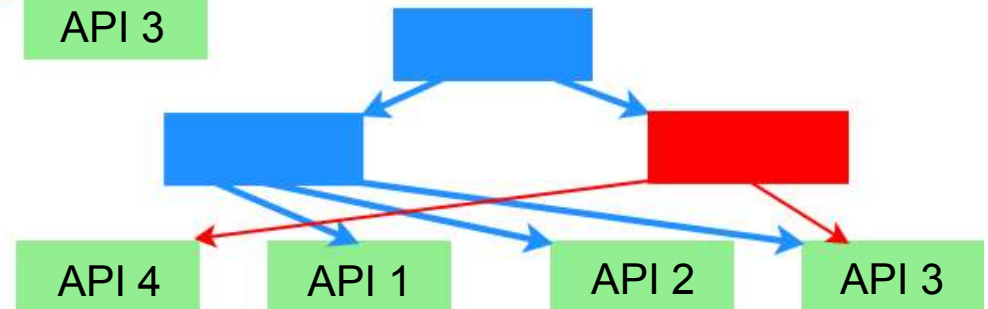
Sample B



# Call Tree Signatures

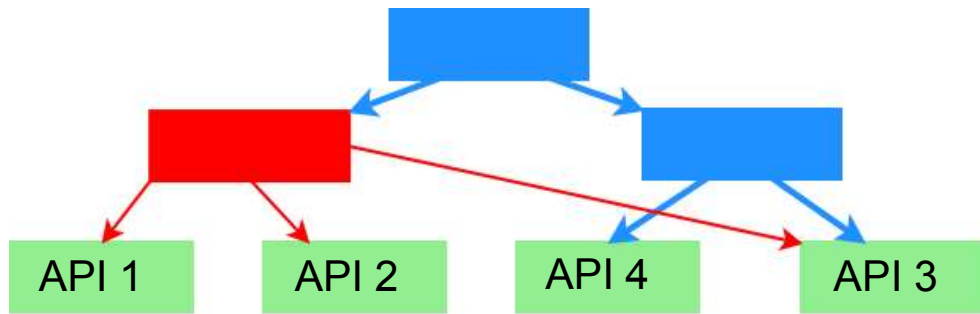


Sample A

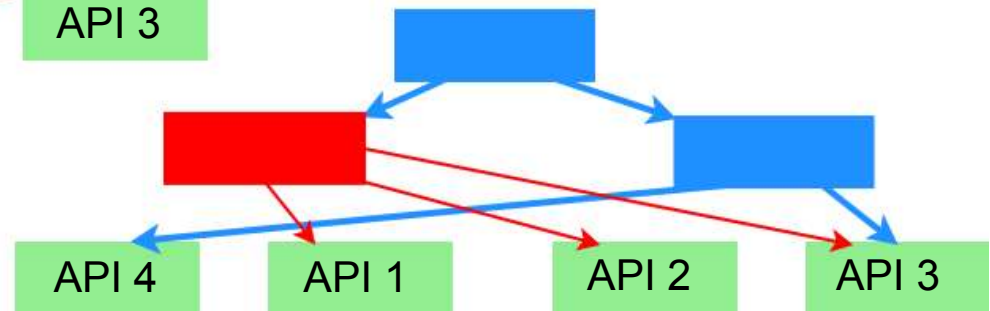


Sample B

# Call Tree Signatures



Sample A



Sample B

# Control Flow Graph Signatures

We number the basic blocks forming a function. The signature is the sorted list of edges linking the blocks.

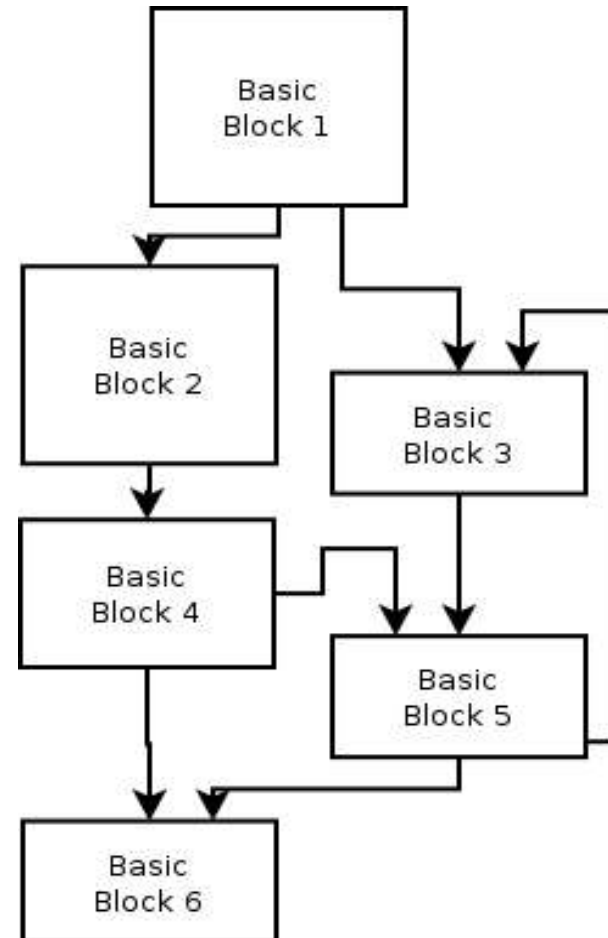
(1,2) (1,3)

(2,4)

(3,5)

(4,5) (4,6)

(5,3) (5,6)



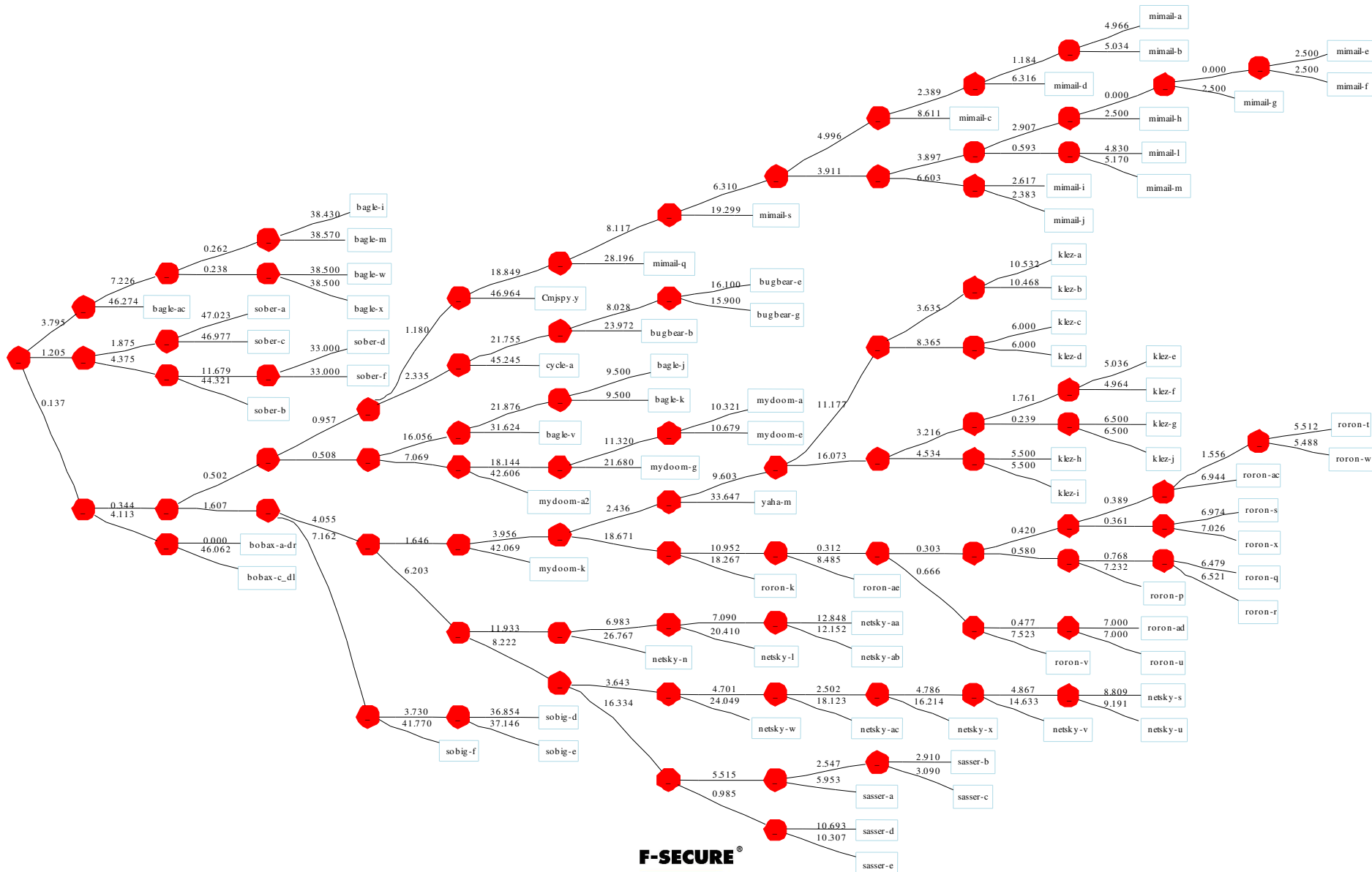
# Possible Uses

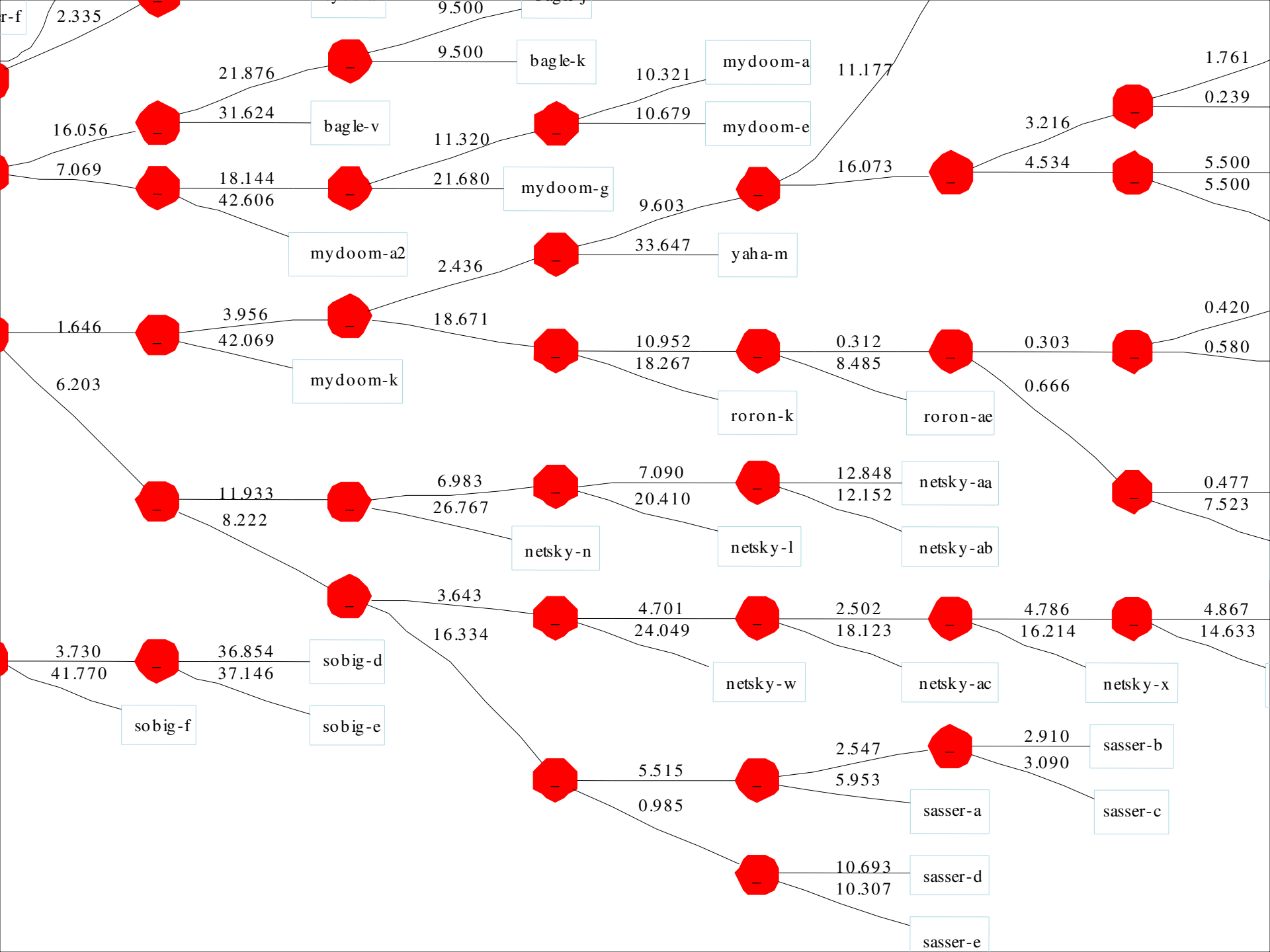
A measure of similarity gives us ways to:

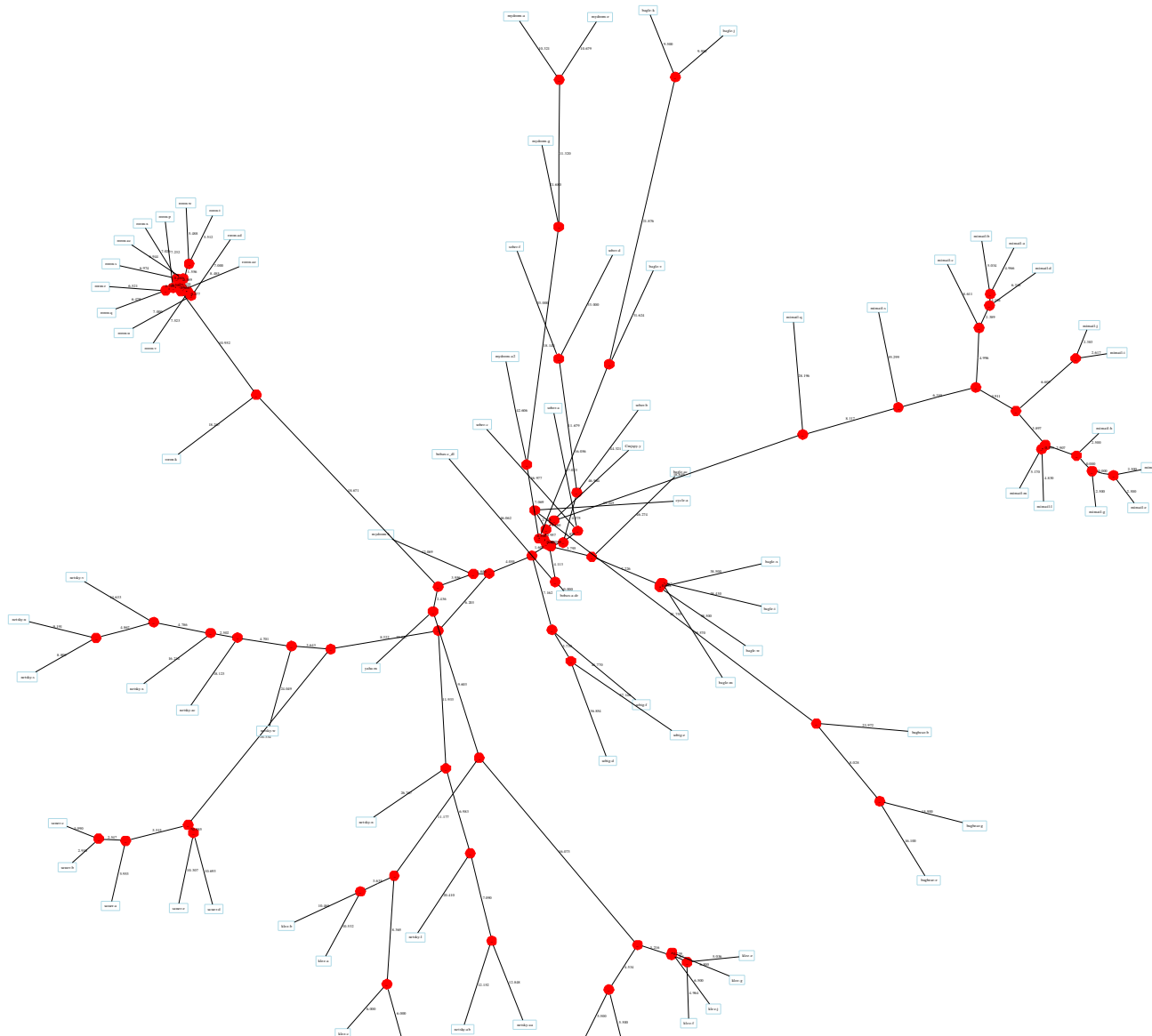
- Classify: Phylogenetics based on similarity matrices
- Identify work by same authors or groups

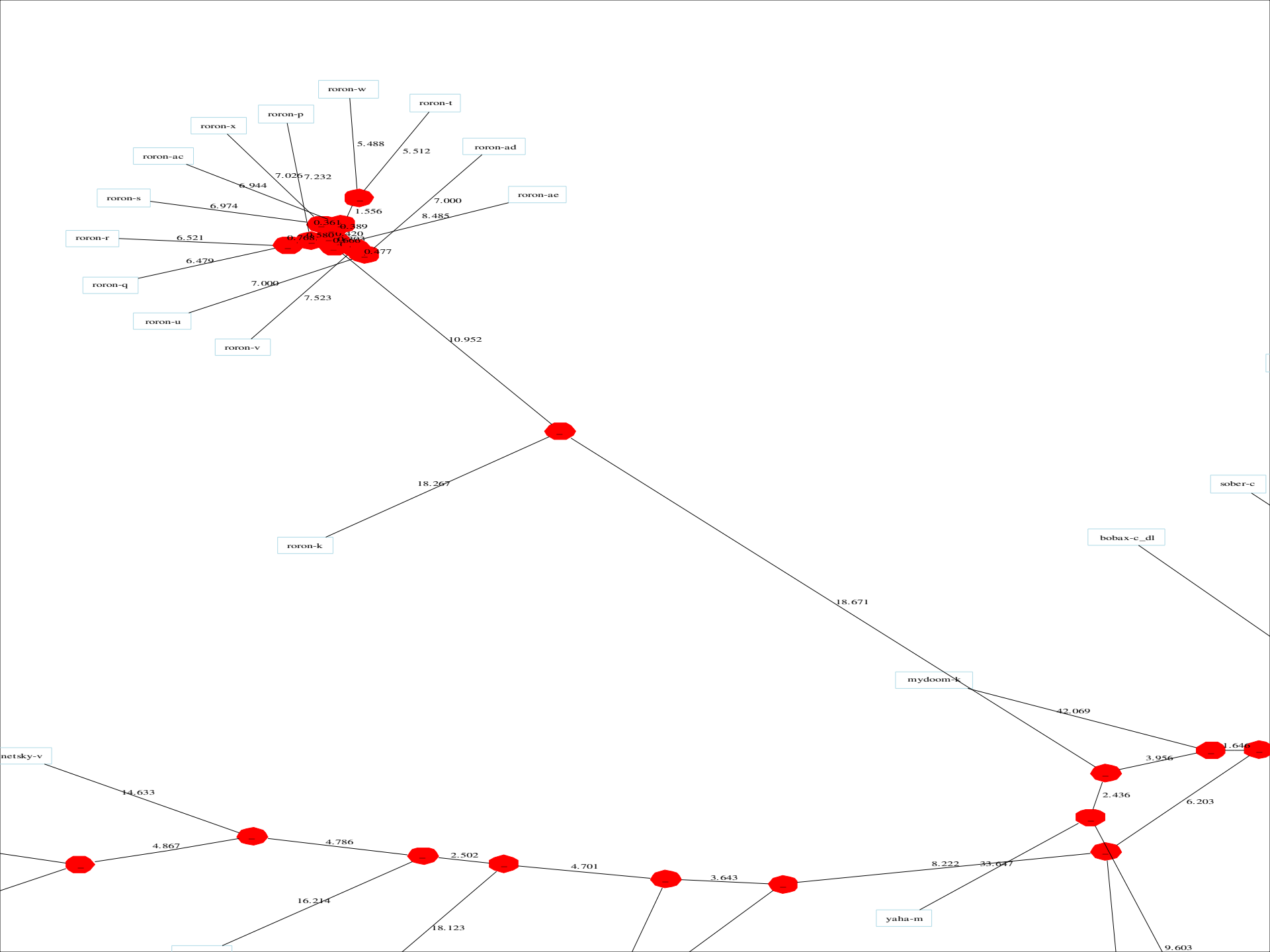
Phylogenetics (From Wikipedia, the free encyclopedia.): *Phylogenetics is the taxonomical classification of organisms based on how closely they are related in terms of evolutionary differences*

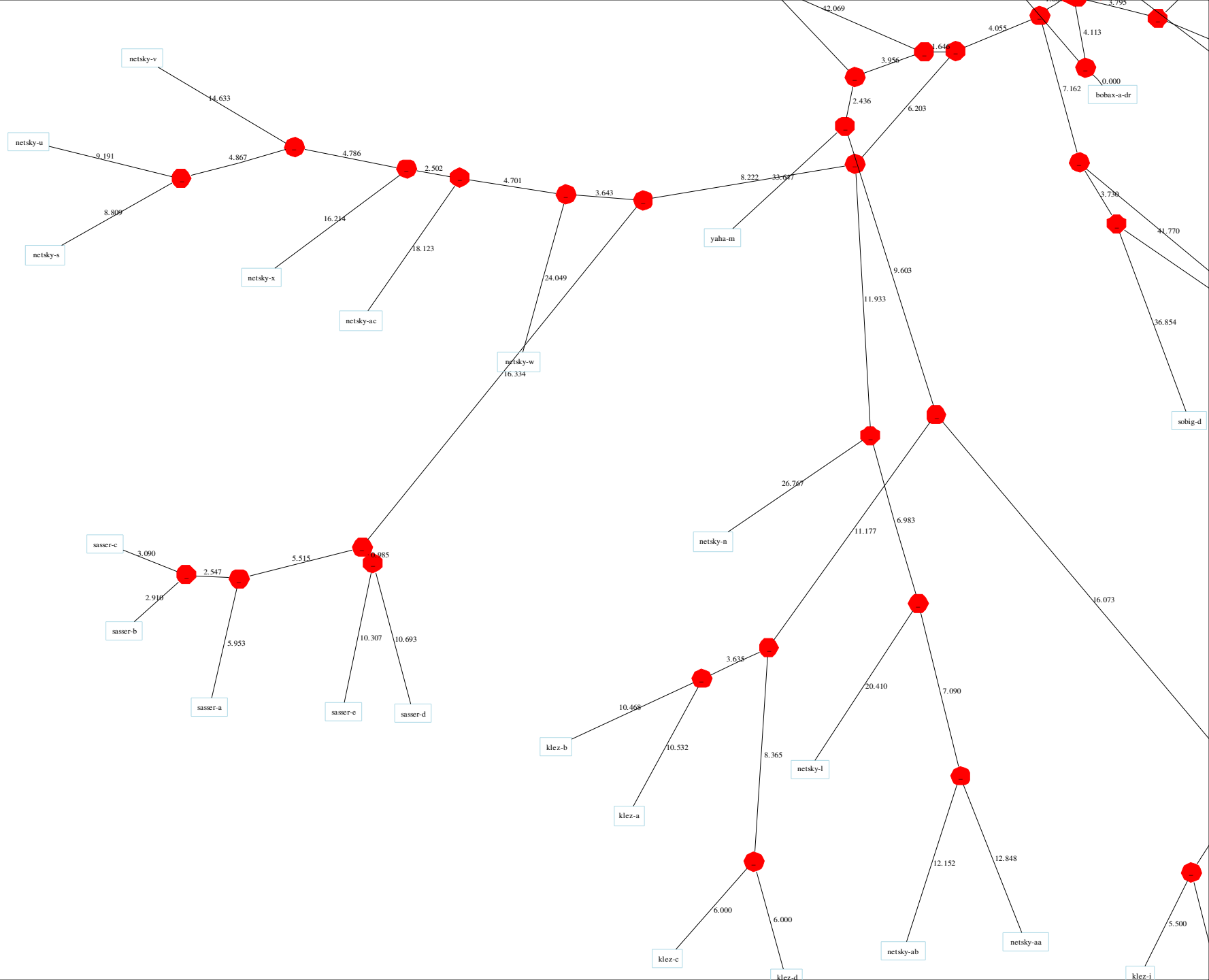


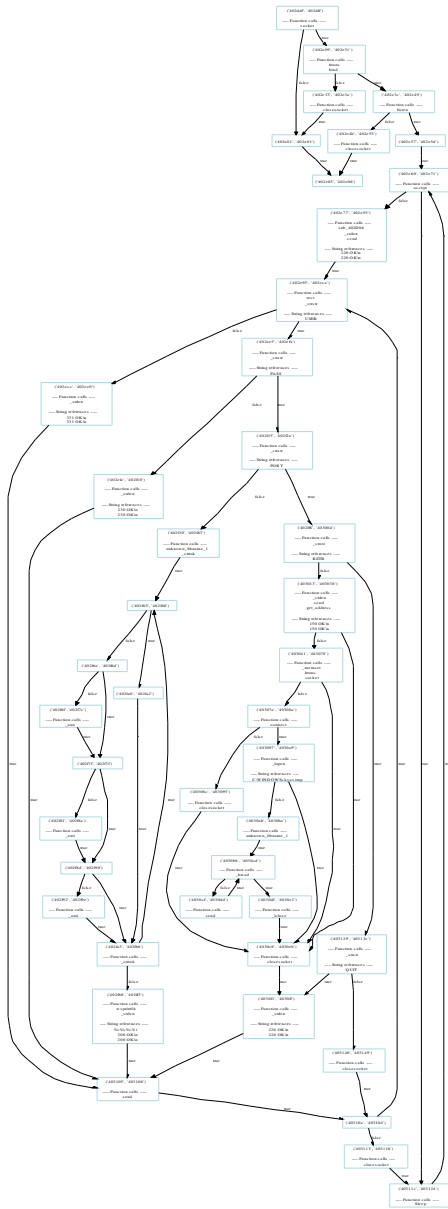




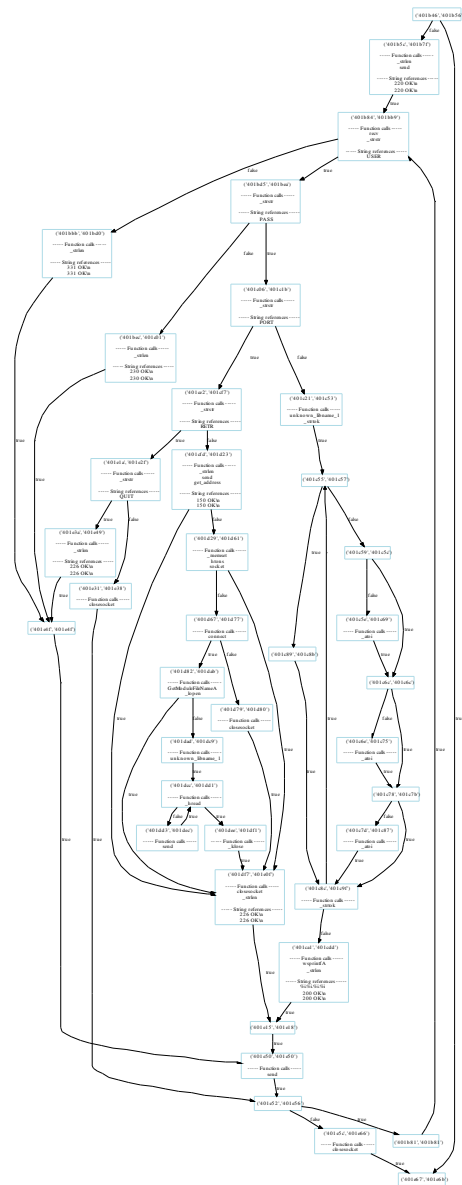








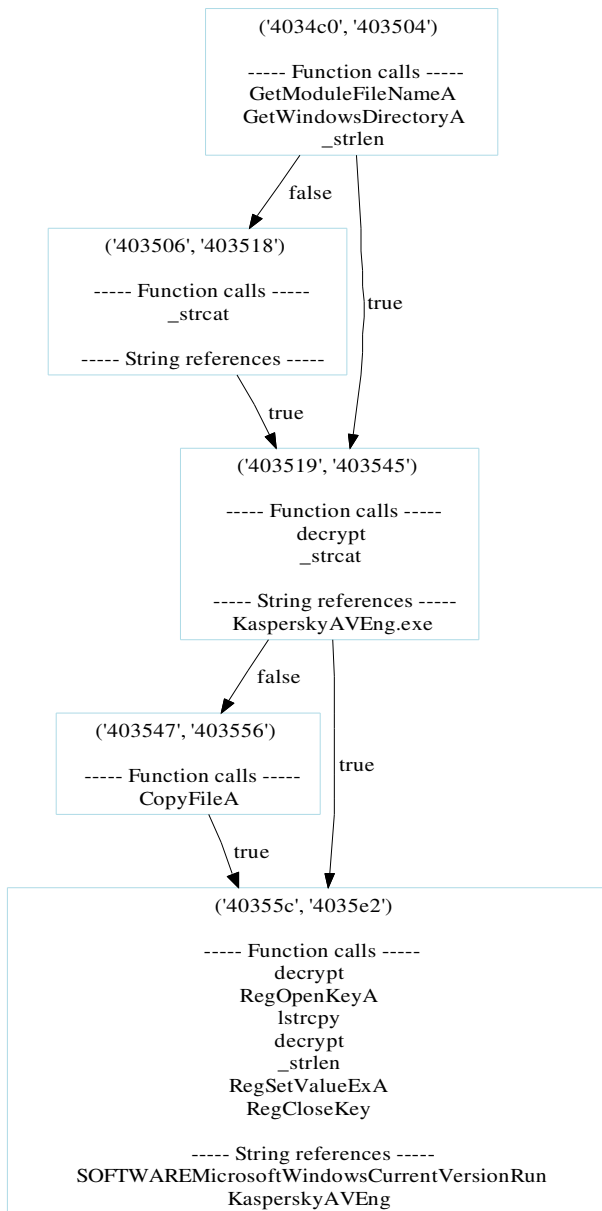
Control Flow Graph for kernel's PTP function



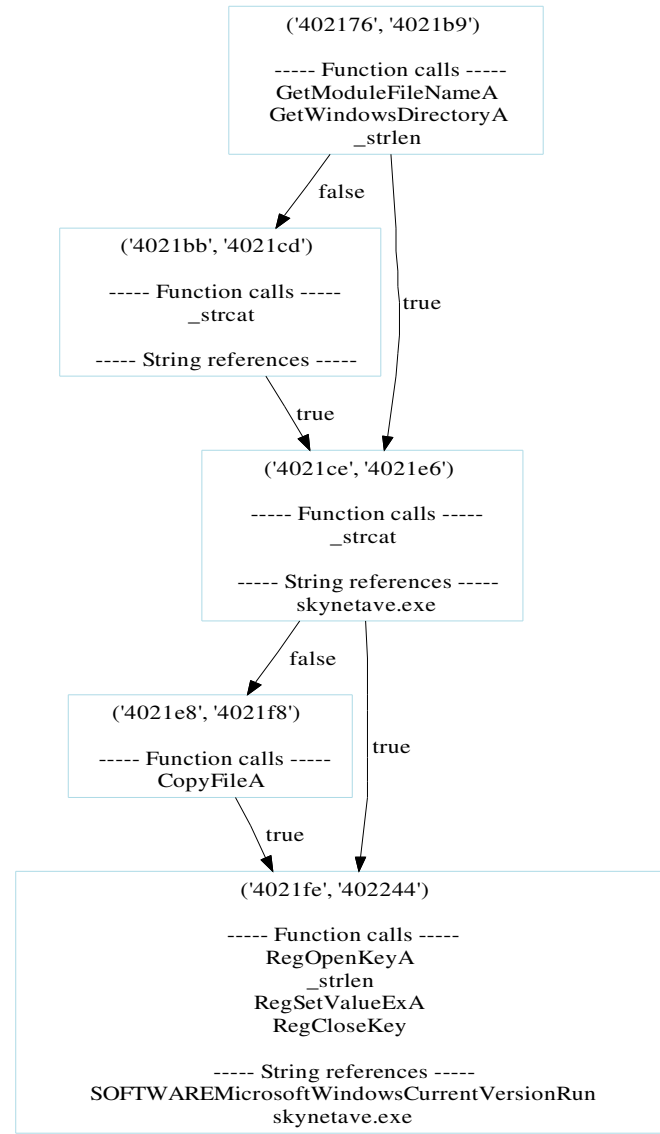
Control Flow Graph for kernel's PTP function







Control Flow Graph for NetSky.V Install function



Control Flow Graph for Sasser.D Install function

# Conclusions

The authors believe approaches such as the one shown have an immense potential

The advantages are patent: compare and classify malware in an automatic manner

Fertile ground of possibilities for taking such methods further



# Closing

IDAPython distribution site:

<http://www.d-dome.net/idapython/>

Questions?

Gergő



Ero

