# BAGUETTE

## HUNTING FOR EVIDENCE OF MALICIOUS BEHAVIORS IN DYNAMIC ANALYSIS REPORTS

VINCENT RAULIN, PIERRE-FRANÇOIS GIMENEZ, YUFEI HAN, VALÉRIE VIET TRIEM TONG
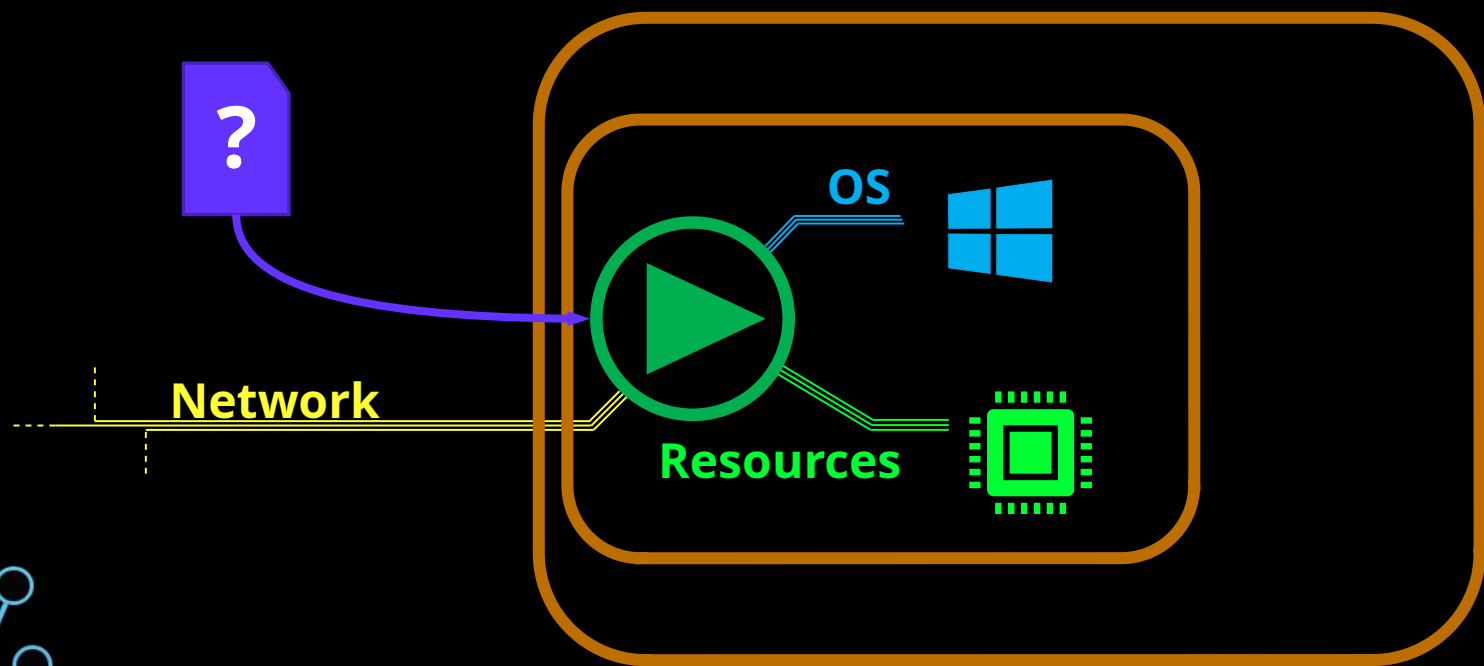
APRIL 5$^{TH}$, THCON 2024

# MALWARE ANALYSIS 101

- >120 million new malware samples per year! (~4/sec) and an estimate of 265 billion USD annually by 2031!

- Exists in many flavors (MS PE, MSI, ELF, JAR archives, Android apps, scripts, PDF, MS Office macros, etc.)

- Two main approaches : static and dynamic analysis

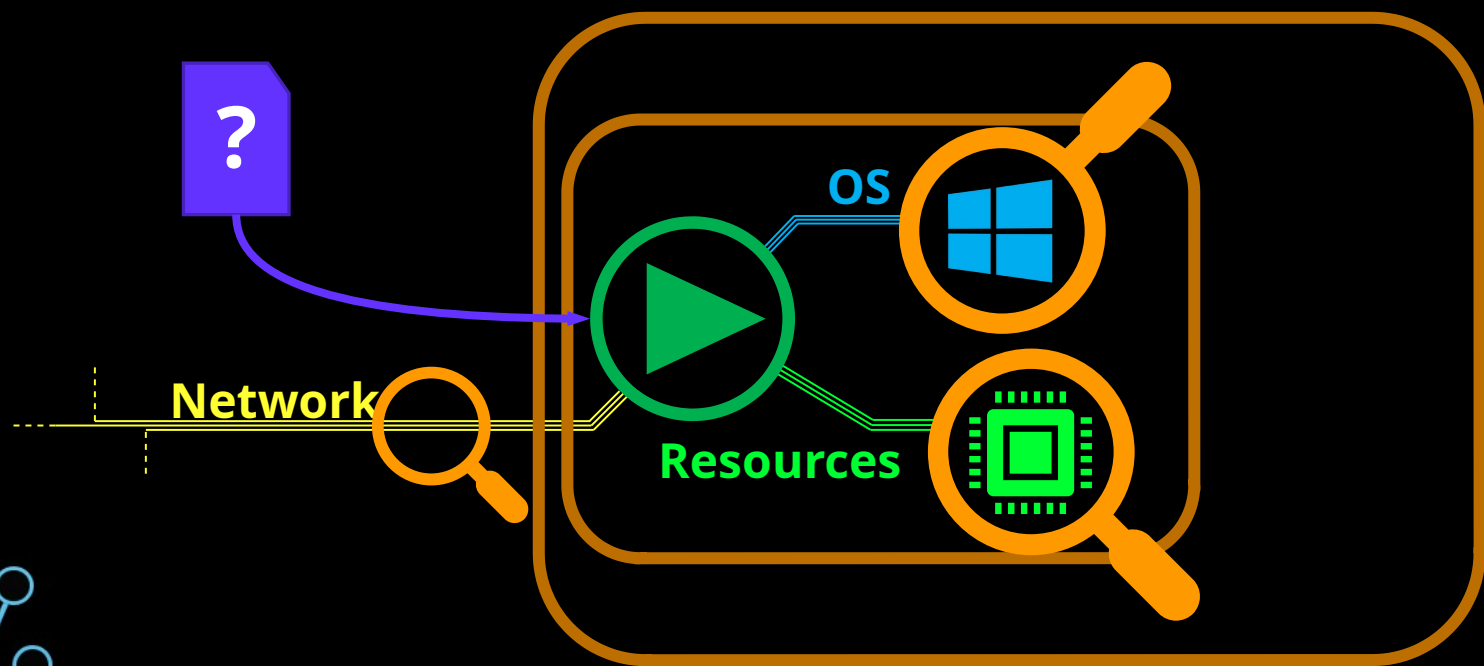- We focus on Windows malware dynamic analysis, using Cuckoo sandbox
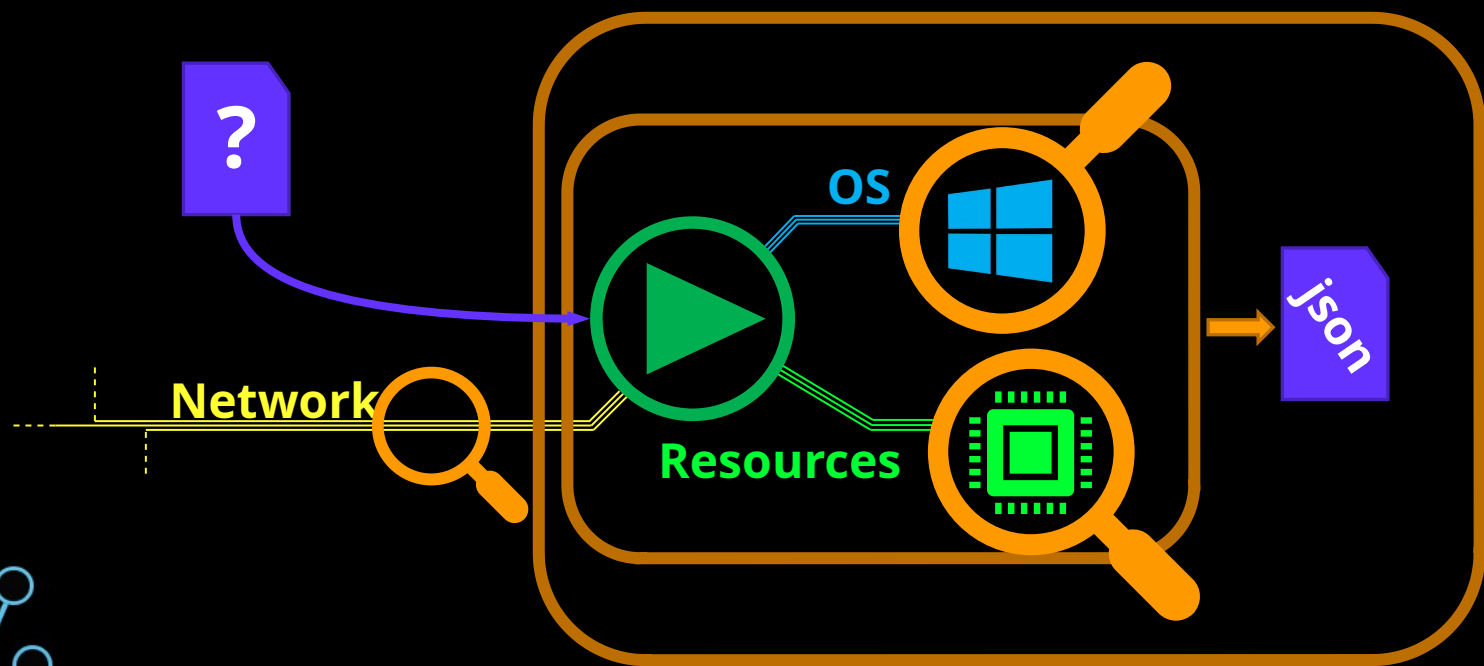
# CUCKOO ANALYSIS PIPELINE

?

# CUCKOO ANALYSIS PIPELINE
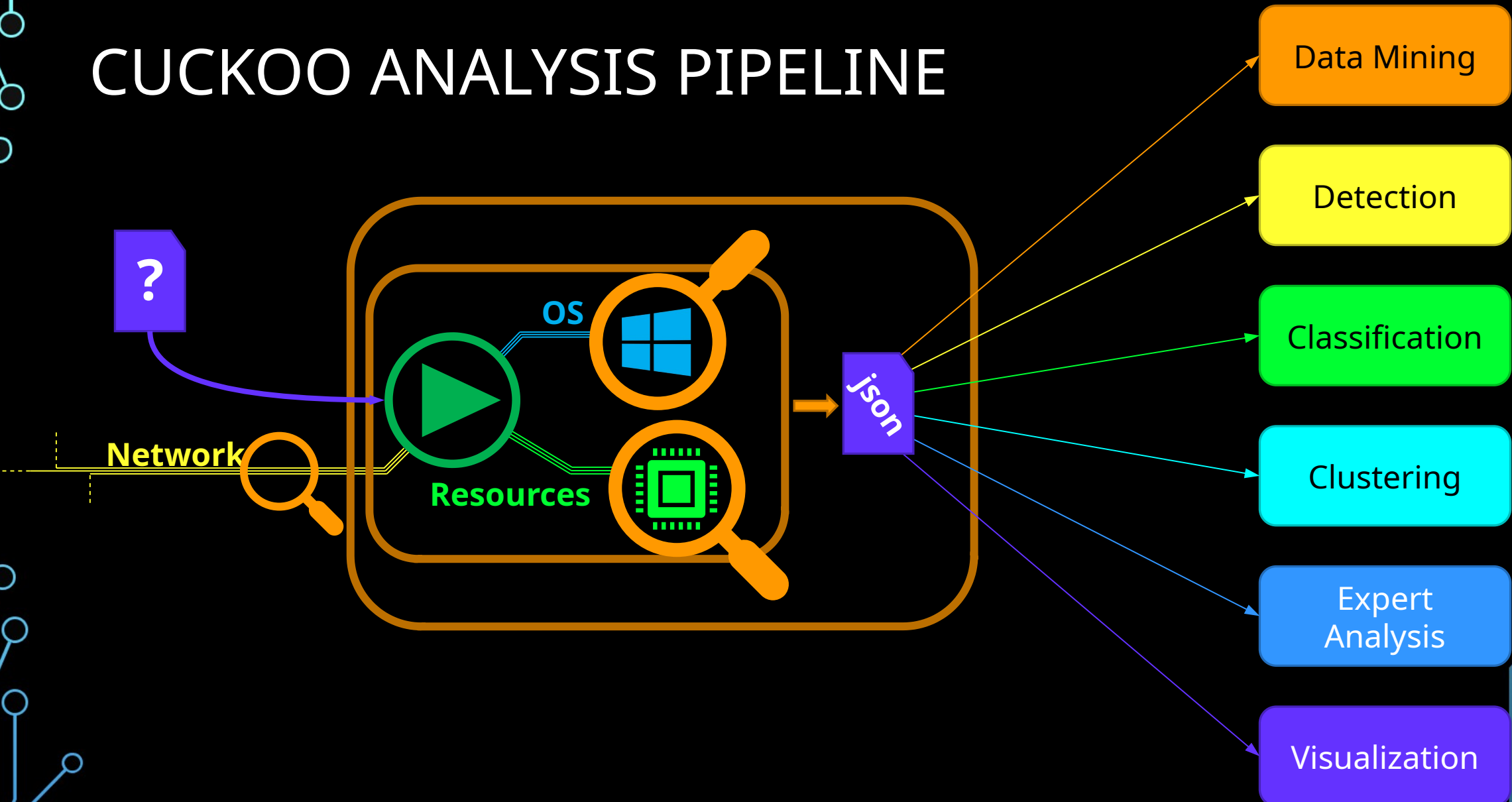
**?**

**OS**

**Network**

**Resources**

# CUCKOO ANALYSIS PIPELINE

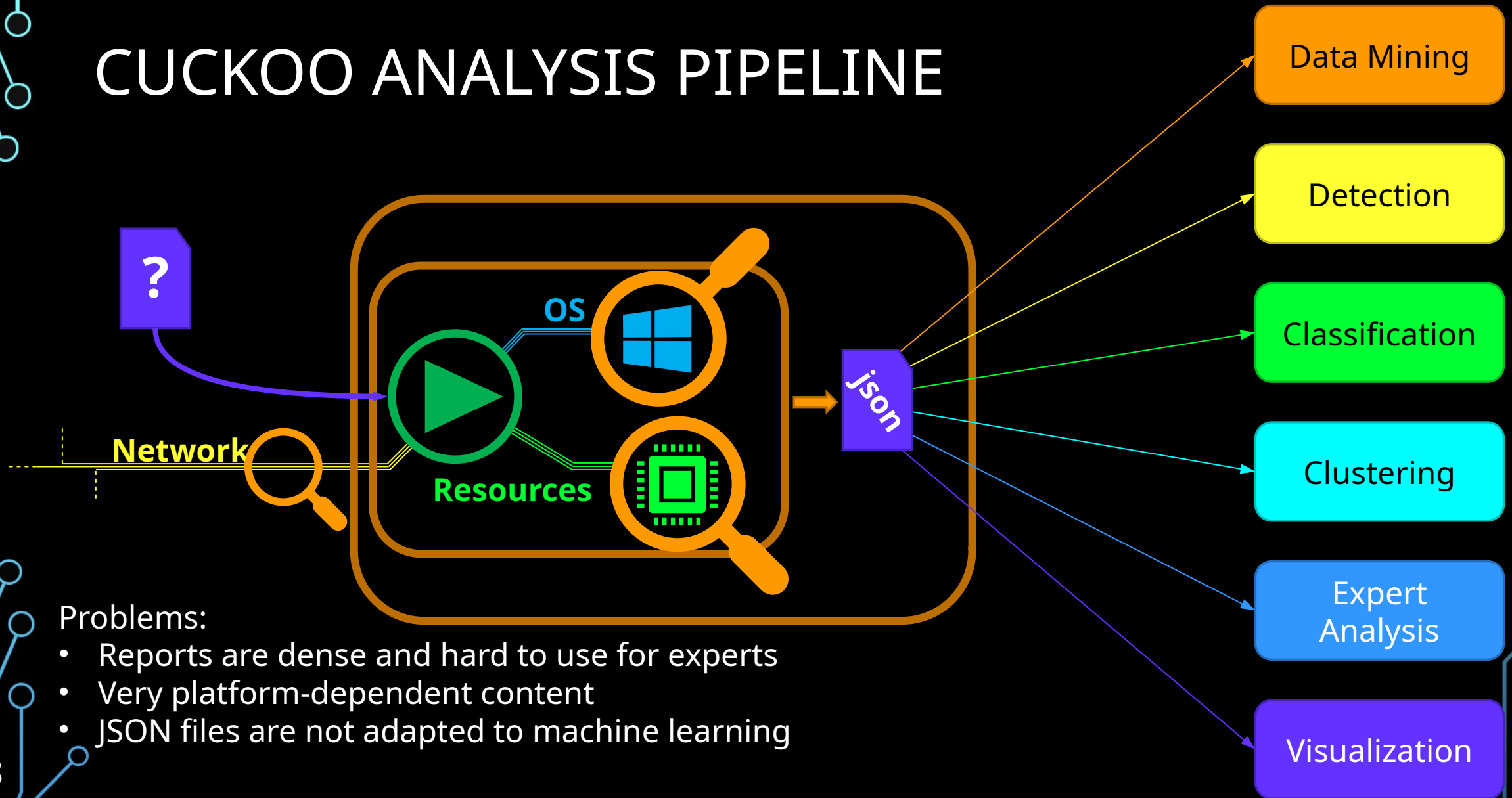# CUCKOO ANALYSIS PIPELINE

CUCKOO ANALYSIS PIPELINE

7

# CUCKOO ANALYSIS PIPELINE



**Problems:**
- Reports are dense and hard to use for experts
- Very platform-dependent content
- JSON files are not adapted to machine learning

8

```json
   1    {
   2  >     "info": { ⋯
  30      },
  31  >     "signatures": [ ⋯
1499      ],
1500  >     "target": { ⋯
1516      },
1517  >     "network": { ⋯
1519      },
1520      "static": {
1521          "pdb_path": null,
1522          "pe_imports": [
1523              {
1524                  "imports": [
1525                      {
1526                          "name": "DeleteCriticalSection",
1527                          "address": "0x40d0b4"
1528                      },
1529                      {
1530                          "name": "LeaveCriticalSection",
1531                          "address": "0x40d0b8"
1532                      },
1533                      {
1534                          "name": "EnterCriticalSection",
1535                          "address": "0x40d0bc"
1536                      },
1537                      {
1538                          "name": "InitializeCriticalSection",
1539                          "address": "0x40d0c0"
1540                      },
1541                      {
1542                          "name": "VirtualFree",
1543                          "address": "0x40d0c4"
1544                      },
1545                      {
1546                          "name": "VirtualAlloc",
1547                          "address": "0x40d0c8"
1548                      },
1549                      {
1550                          "name": "LocalFree",
1551                          "address": "0x40d0cc"
1552                      },
1553                      {
1554                          "name": "LocalAlloc",
1555                          "address": "0x40d0d0"
```

```json
1500 >        "target": {…
1516          },
1517 >        "network": {…
1519          },
1520 v        "static": {
1521              "pdb_path": null,
1522 >            "pe_imports": […
1947              ],
1948              "peid_signatures": null,
1949              "keys": [],
1950              "signature": [],
1951              "pe_timestamp": "1992-06-20 00:22:17",
1952              "pe_exports": [],
1953              "imported_dll_count": 8,
1954              "pe_imphash": "884310b1928934402ea6fec1dbd3cf5e",
1955 >            "pe_resources": […
2044              ],
2045 >            "pe_versioninfo": […
2078              ],
2079 v            "pe_sections": [
2080                  {
2081                      "size_of_data": "0x00009400",
2082                      "virtual_address": "0x00001000",
2083                      "entropy": 6.557291120606633,
2084                      "name": "CODE",
2085                      "virtual_size": "0x0000933c"
2086                  },
2087 v                {
2088                      "size_of_data": "0x00000400",
2089                      "virtual_address": "0x0000b000",
2090                      "entropy": 2.7679914923058857,
2091                      "name": "DATA",
2092                      "virtual_size": "0x0000024c"
2093                  },
2094 v                {
2095                      "size_of_data": "0x00000000",
2096                      "virtual_address": "0x0000c000",
2097                      "entropy": 0.0,
2098                      "name": "BSS",
2099                      "virtual_size": "0x00000e4c"
2100                  },
2101 v                {
2102                      "size_of_data": "0x00000a00",
2103                      "virtual_address": "0x0000d000",
2104                      "entropy": 4.430733069799032,
```
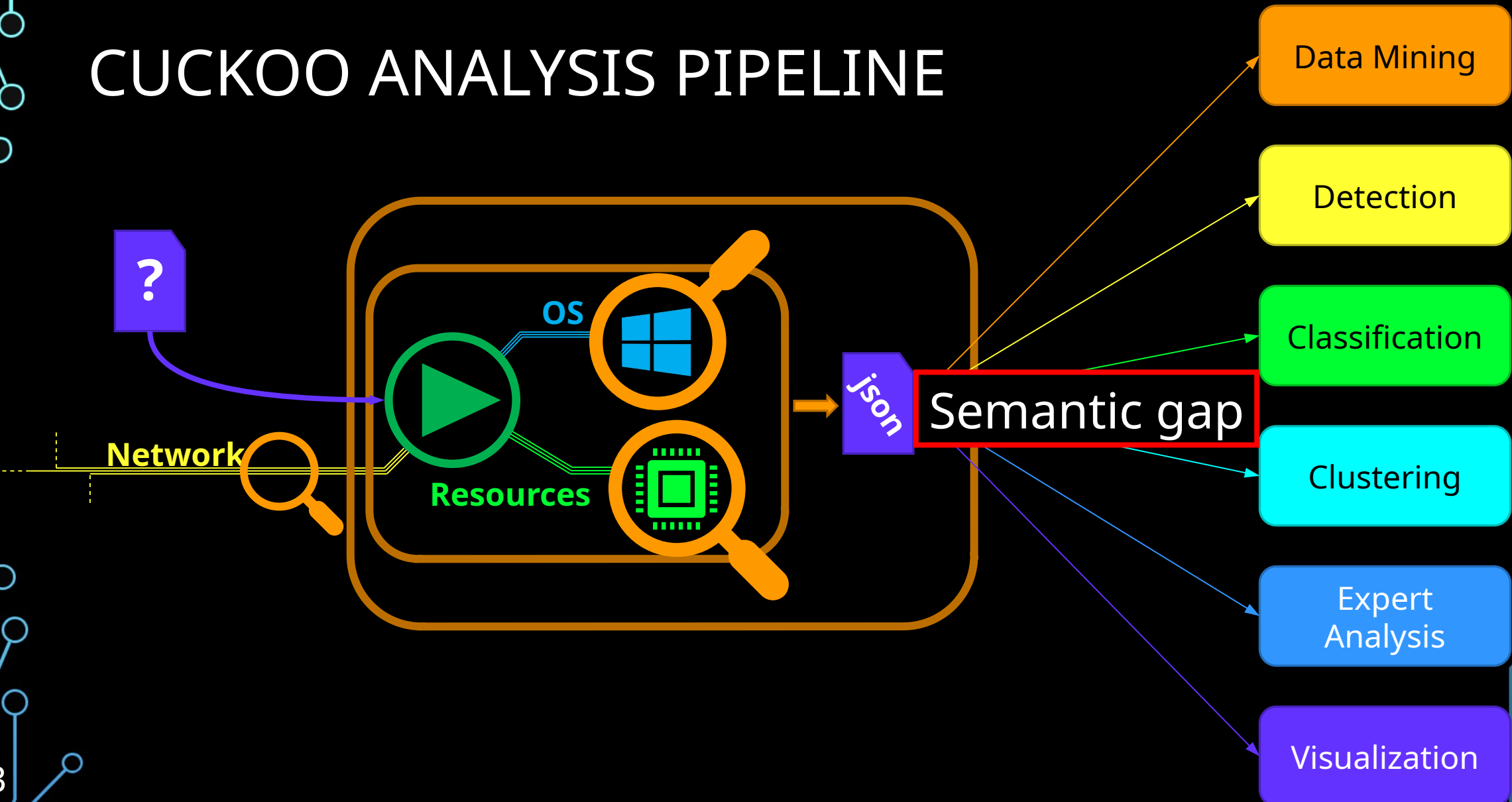
10

```json
{
    "info": {…
    },
    "signatures": […
    ],
    "target": {…
    },
    "network": {…
    },
    "static": {…
    },
    "dropped": […
    ],
    "behavior": {
        "generic": [
            {
                "process_path": "C:\\Users\\Marc Elbichon\\AppData\\Local\\Temp\\is-6BR7Q.tmp\\c557b1737ec0d359ebc4868caf8a31c2a31ec56954ba51bb39d90312d023dc97.tmp",
                "process_name": "c557b1737ec0d359ebc4868caf8a31c2a31ec56954ba51bb39d90312d023dc97.tmp",
                "pid": 2968,
                "summary": {
                    "file_created": [
                        "C:\\Users\\Marc Elbichon\\AppData\\Local\\Temp\\is-S7NDS.tmp\\_isetup\\_shfoldr.dll",
                        "C:\\Users\\Marc Elbichon\\AppData\\Local\\Temp\\is-S7NDS.tmp\\_isetup\\_setup64.tmp",
                        "C:\\Users\\Marc Elbichon\\AppData\\Local\\Temp\\is-S7NDS.tmp\\_isetup\\_RegDLL.tmp",
                        "C:\\Program Files (x86)\\FevsoftFR\\FinalRecovery\\data\\is-5CVH4.tmp",
                        "C:\\Program Files (x86)\\FevsoftFR\\FinalRecovery\\is-S2PNI.tmp",
                        "C:\\Program Files (x86)\\FevsoftFR\\FinalRecovery\\is-9LRAI.tmp",
                        "C:\\Users\\Marc Elbichon\\AppData\\Local\\Temp\\is-S7NDS.tmp\\_isetup\\_iscrypt.dll",
                        "C:\\Program Files (x86)\\FevsoftFR\\FinalRecovery\\unins000.dat",
                        "C:\\Program Files (x86)\\FevsoftFR\\FinalRecovery\\is-V1A4O.tmp",
                        "C:\\Program Files (x86)\\FevsoftFR\\FinalRecovery\\is-MQUFU.tmp",
                        "C:\\Program Files (x86)\\FevsoftFR\\FinalRecovery\\is-5F2AK.tmp"
                    ],
                    "file_recreated": [
                        "C:\\Program Files (x86)\\FevsoftFR\\FinalRecovery\\unins000.dat",
                        "\\Device\\KsecDD",
                        "\\Device\\DeviceApi\\CMApi"
                    ],
                    "directory_created": [
                        "C:\\Users\\Marc Elbichon\\AppData\\Local\\Temp\\is-S7NDS.tmp",
                        "C:\\Program Files (x86)\\FevsoftFR",
                        "C:\\Program Files (x86)\\FevsoftFR\\FinalRecovery\\data",
                        "C:\\Users\\Marc Elbichon\\AppData\\Local\\Temp\\is-S7NDS.tmp\\_isetup",
                        "C:\\Program Files (x86)\\FevsoftFR\\FinalRecovery"
                    ],
```
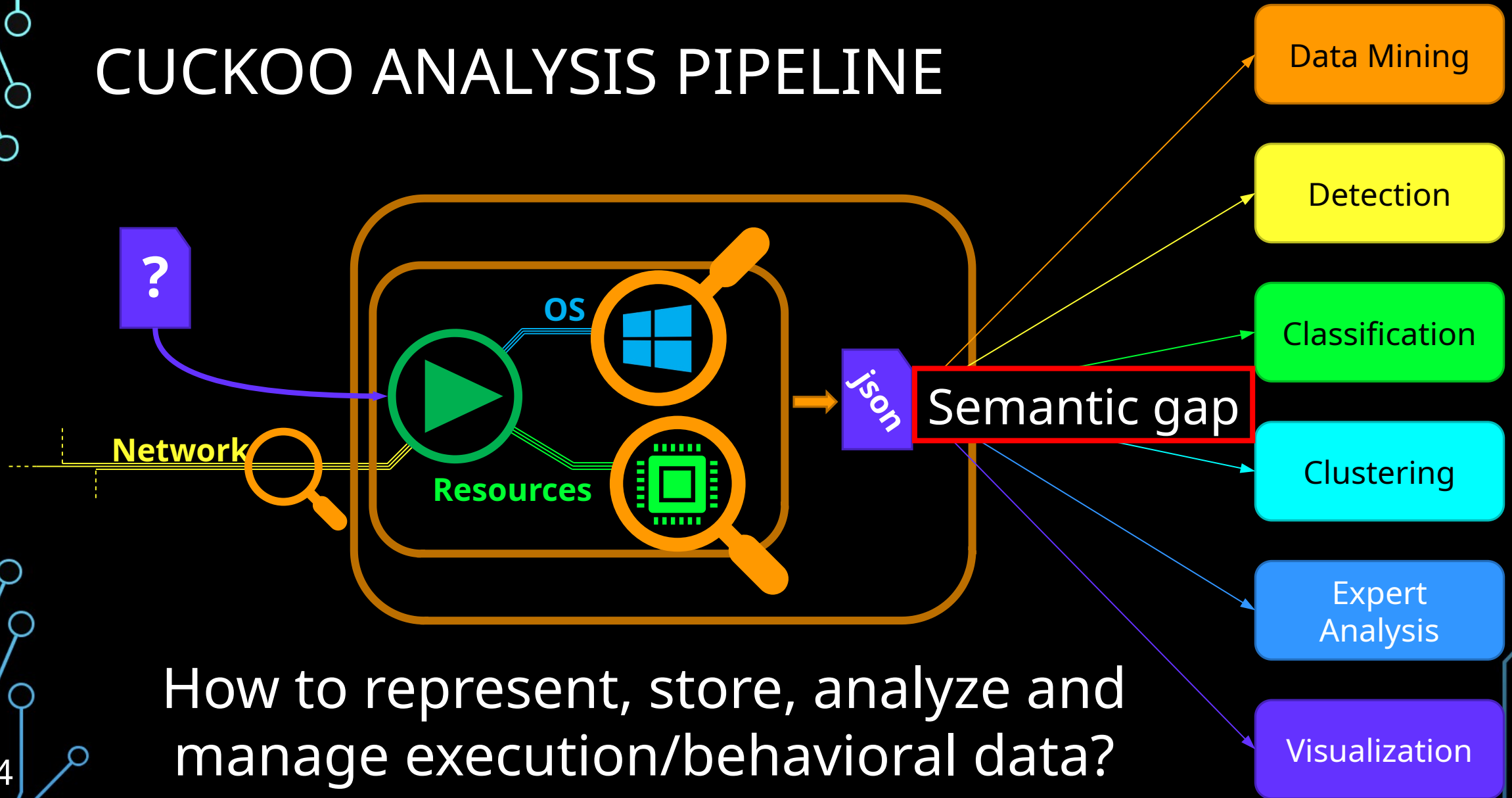
11

```
3867            {
3868                "process_path": "C:\\Users\\Marc Elbichon\\AppData\\Local\\Temp\\c557b1737ec0d359ebc4868caf8a31c2a31ec56954ba51bb39d90312d023dc97.exe",
3869                "calls": [
3870                    {
3871                        "category": "system",
3872                        "status": 1,
3873                        "stacktrace": [],
3874                        "api": "LdrGetDllHandle",
3875                        "return_value": 0,
3876                        "arguments": {
3877                            "module_name": "kernel32.dll",
3878                            "stack_pivoted": 0,
3879                            "module_address": "0x75d20000"
3880                        },
3881                        "time": 1674405919.603929,
3882                        "tid": 6644,
3883                        "flags": {}
3884                    },
3885                    {
3886                        "category": "process",
3887                        "status": 1,
3888                        "stacktrace": [],
3889                        "api": "NtAllocateVirtualMemory",
3890                        "return_value": 0,
3891                        "arguments": {
3892                            "process_identifier": 5356,
3893                            "region_size": 1048576,
3894                            "stack_dep_bypass": 0,
3895                            "stack_pivoted": 0,
3896                            "heap_dep_bypass": 0,
3897                            "protection": 1,
3898                            "process_handle": "0xffffffff",
3899                            "allocation_type": 8192,
3900                            "base_address": "0x00e70000"
3901                        },
3902                        "time": 1674405919.603929,
3903                        "tid": 6644,
3904                        "flags": {
3905                            "protection": "PAGE_NOACCESS",
3906                            "allocation_type": "MEM_RESERVE"
3907                        }
3908                    },
3909                    {
3910                        "category": "process",
3911                        "status": 1,
```

It is difficult to map these data to the actual effect the software has on the system
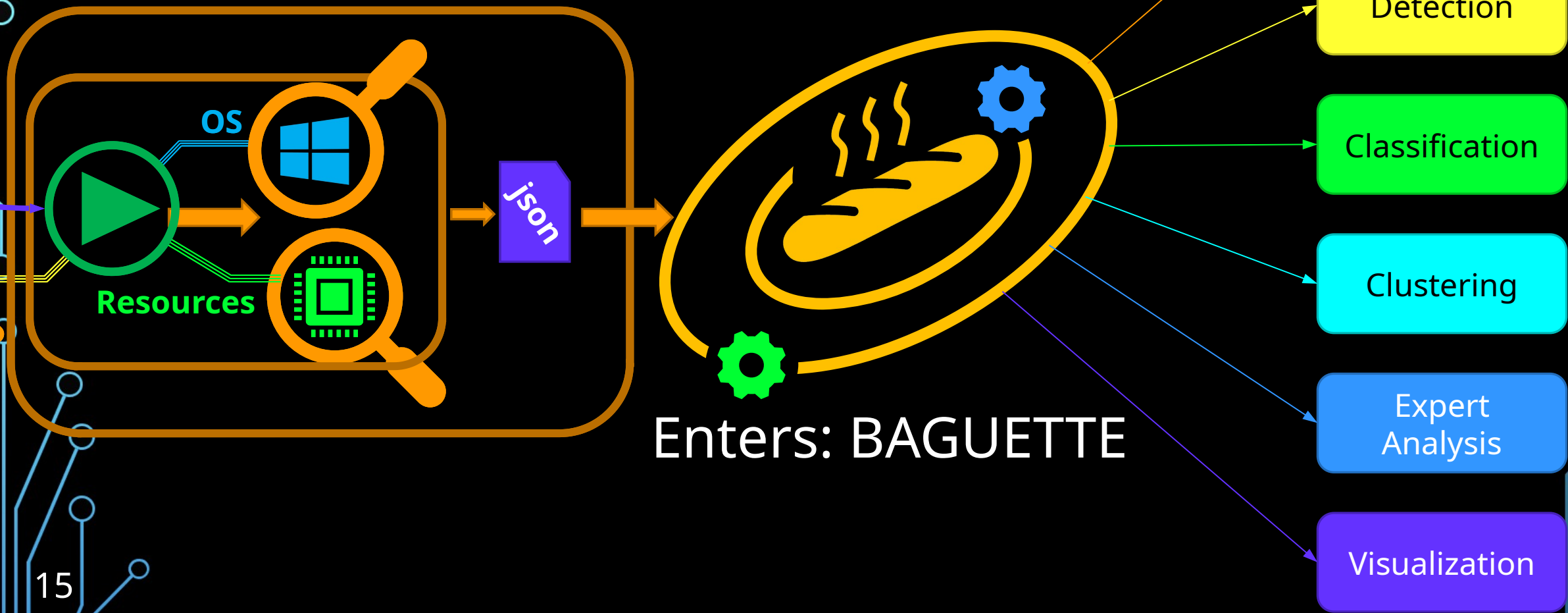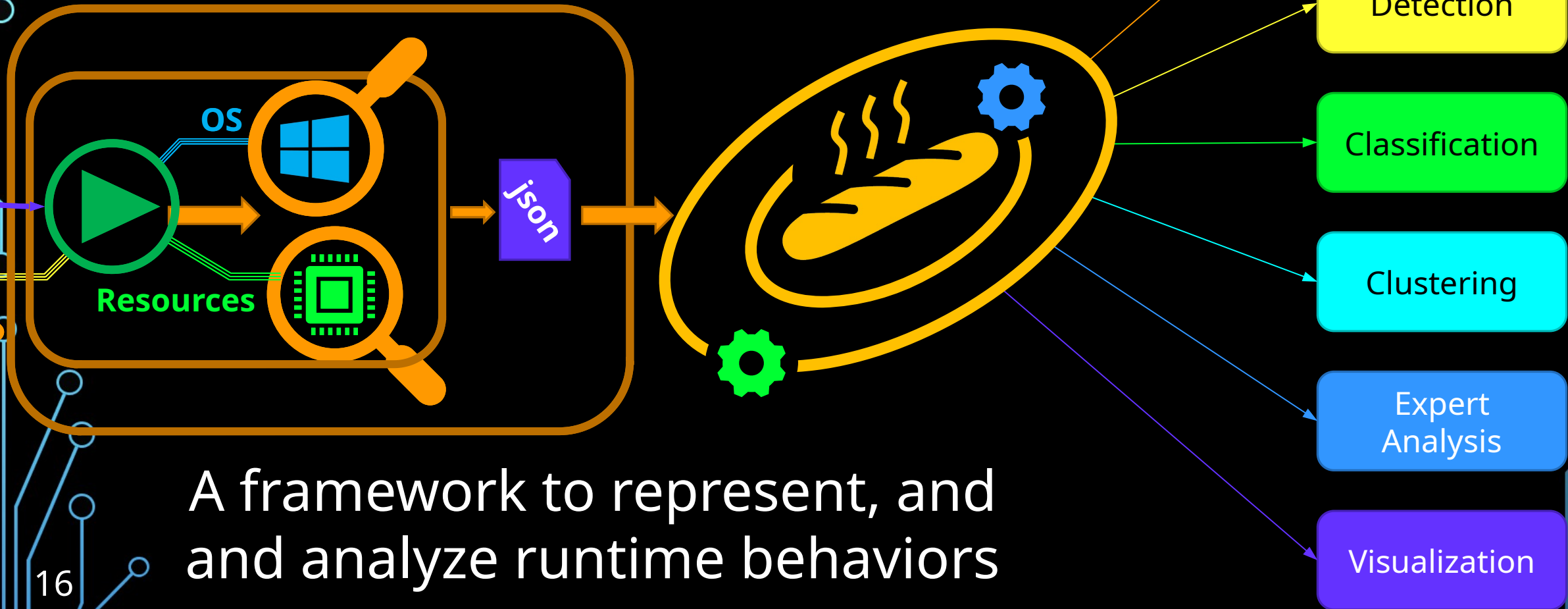
# CUCKOO ANALYSIS PIPELINE

13

CUCKOO ANALYSIS PIPELINE

How to represent, store, analyze and manage execution/behavioral data?

14

CUCKOO ANALYSIS PIPELINE

OS

Resources
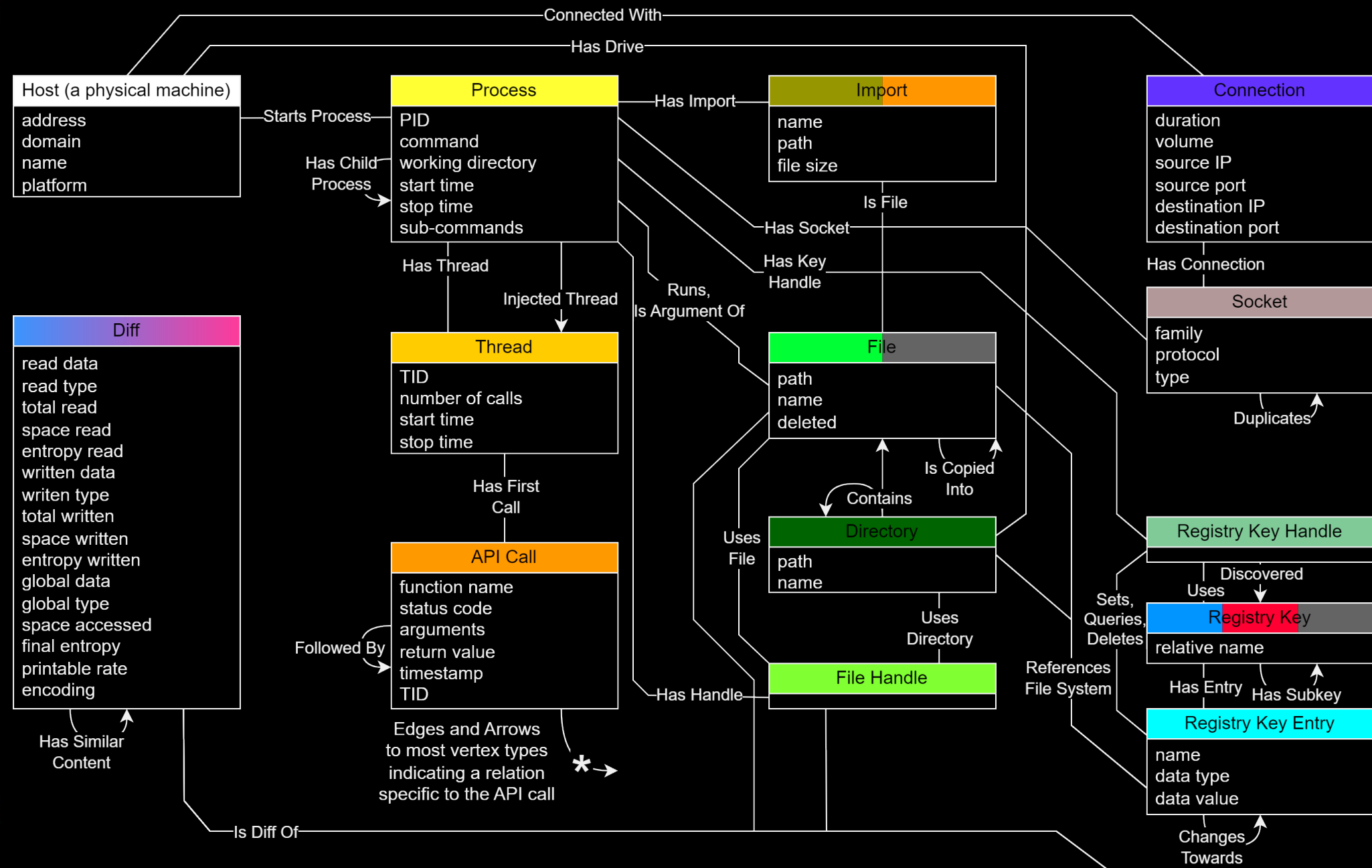
json

Enters: BAGUETTE

Data Mining

Detection

Classification

Clustering

Expert Analysis

Visualization

15

# CUCKOO ANALYSIS PIPELINE

A framework to represent, and and analyze runtime behaviors

16

Behavioral Analysis Graph Using Execution Traces Towards Explainability

# WHAT IS BAGUETTE?

- A graphical representation of dynamic analysis traces

- Heterogeneous graph

- Shows resources given by the OS: file system, registry keys, network connection, etc.

- Links related resources

# DIFF NODES

- Any high-level read or write operations involve many elementary read and write system calls

- We merge them into « diff » nodes summarizing data transfer

- Diff nodes include all the data read from and written into a socket or a file

- That way, we can easily analyzed read and written data

  - Entropy computations

  - ASCII or binary data

  - Header analysis

# BAGUETTE EXAMPLE

# METAGRAPHS TO ANALYZE A BAGUETTE

- We use graph patterns called « metagraphs »
  - They are graphs where:
    - nodes can match one or several BAGUETTE nodes
    - edges can match one or several BAGUETTE edges
    - nodes can also have conditions of BAGUETTE attributes
- Since BAGUETTE are high-level, we can manually write metagraphs that match specific behaviors

# METAGRAPH EXAMPLES



High Entropy Writing

Vertex[Diff]

Edge[builder.source.types.data.IsDiffOf | builder.source.types.data.WritesInto]

Vertex[Diff]{lambda x : x.written_entropy >= 6.0 and x.written_space > 0}

Edge[builder.source.types.data.IsDiffOf | builder.source.types.data.WritesInto]

Vertex[File]

Vertex[Handle]

Edge[HasHandle]

Vertex[Key]{lambda x : x.parent_key is None and x.name.upper() in {'HKEY_LOCAL_MACHINE', 'HKEY_CURRENT_USER'}}

Vertex[Key]{lambda x : x.name.lower() == 'microsoft'}

Edge[UsesFile]

Vertex[Process]

Arrow[HasSubKey]

Arrow[HasSubKey]

Vertex[Process]

Vertex[Key]{lambda x : x.name.lower() == 'windows'}

Arrow[HasSubKey]

Arrow[HasSubKey]

Edge[Runs]

Auto-Run

Vertex[Key]{lambda x : x.name.lower() == 'software'}

Vertex[File]

Arrow[HasChildProcess]

Edge[Runs]

Vertex[Process]

Arrow[HasSubKey]

Vertex[Key]{lambda x : x.name.lower() == 'currentversion'}

Vertex[File]{lambda x : x.extension not in {'exe', 'vbs', 'ps1'}}

Extraction and Execution

Covert Execution

Arrow[HasSubKey]

Vertex[Key]{lambda x : x.name.lower() in {'run', 'runonce'}}

Edge[HasEntry]

Vertex[KeyEntry]

Vertex[File]

Edge[SetsEntry]

Edge[ReferencesFileSystem]

Edge[IsDiffOf]

Changed File Type

Vertex[Handle]

Vertex[File]

Vertex[Diff]{lambda x : x.read_type != x.written_type}

22

# EXPERIMENTS

- We analyze three malware families:
  - GCleaner, a file dropper
  - SnakeKeyLogger, a key logger and spyware
  - LockBit, a ransomware

| Metagraph | GCleaner (247) | | | SnakeKeyLogger (436) | | | LockBit (7) | | |
|---|---|---|---|---|---|---|---|---|---|
| | p | n | σ | p | n | σ | p | n | σ |
| High-Entropy Writing | 97.57% | 1.53 | 0.59 | 13.76% | 1.08 | 0.28 | 28.57% | 2450.0 | 1878.0 |
| Changed File Type | 97.57% | 1.0 | 0.0 | 4.82% | 1.05 | 0.21 | 14.29% | 1.0 | 0.0 |
| Covert Execution | 98.38% | 1.0 | 0.0 | 0% | - | - | 0% | - | - |
| Extraction and Execution | 98.38% | 2.97 | 0.17 | 13.53% | 1.0 | 0.0 | 0% | - | - |
| Auto-Run | 0% | - | - | 0% | - | - | 28.57% | 1.0 | 0.0 |
| p : Proportion of matches, n : average number per matching sample, σ : standard deviation per matching sample | | | | | | | | | |

- Quite different proportions depending on families

- Tells us how to select samples (for example, which sample executed their payloads)

25

# WHAT NEXT?

This was the state of our research at the time we submitted to THCon... But since, we had some fun

Next research question: how to automatically create metagraphs from a dataset of malware ?

Still a work in progress

# NEW GOAL

Given a unlabelled dataset of malware samples, how can we:

- Recover clusters of behaviors that hopefully match the families/classes
- Recover the behavioral patterns characteristic of each of these clusters

With some constrains:

- In an unsupervised way (no labels)
- Without expert knowledge on malware analysis, just on system programming

# EXPERIMENTAL DETAILS

- Dataset: 13 families, 100 samples for each family
- MetaGraph library generation is an iterative process:
  - Generate new valid metagraphs from the previous library
  - Search them across the BAGUETTE dataset
  - Select the best ones according to some metrics
  - Repeat
- A classic genetic algorithm, with two hurdles:
  - How to mutate metagraphs?
  - How to select the best ones?

# EXPERIMENTAL DETAILS

- What is a good metagraph?
  - A metagraph that is rarely present? → probably not significative
  - A metagraph that is always present? → probably typical Windows behavior, like DLL imports, etc., not interesting for malware analysis
  - A metagraph that are very common in some software but very rare in others? → sounds like something akin to a signature!

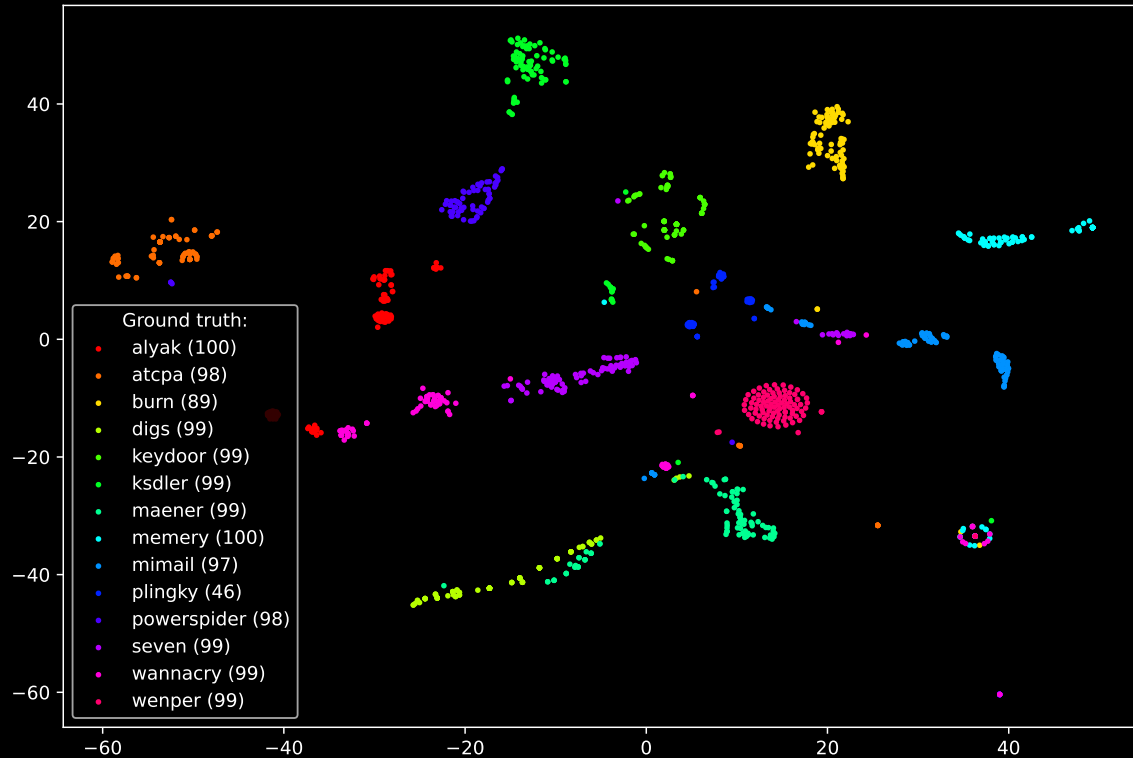- But what does it mean, mathematically?

# EXPERIMENTAL DETAILS

- We experimented with several metrics, I'll describe the best one
- TF-ISF (Term-Frequency / Inverse Sample Frequency) inspired from TF-IDF
- For one metagraph, it's the multiplication of two terms:
  - The number of occurrences in all BAGUETTEs
  - The logarithm of the inverse of the number of BAGUETTEs matching this metagraph
- Intuitively, we want common metagraphs that are only present in a few BAGUETTEs
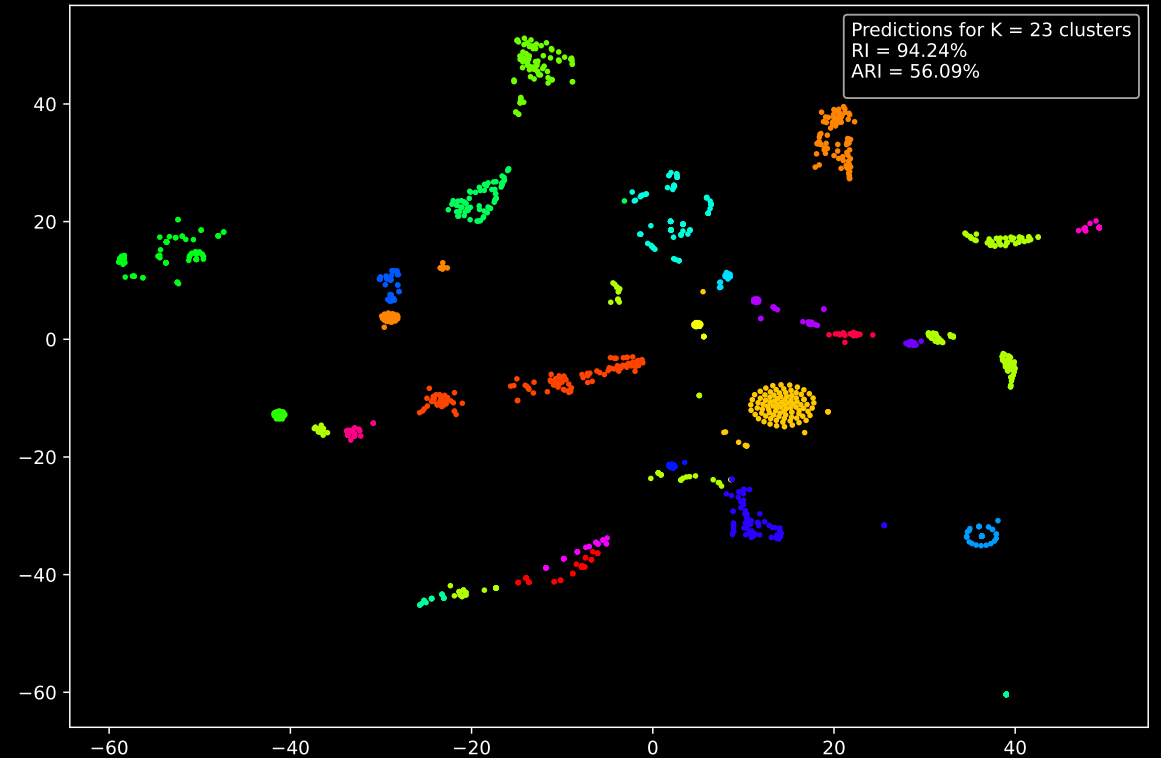
# EXPERIMENTAL DETAILS

- Clustering is made in a vector space where each BAGUETTE graph is described by the number of matches for each metagraph

- So, if we have 100 metagraphs, each BAGUETTE is represented by a vector of 100 numbers

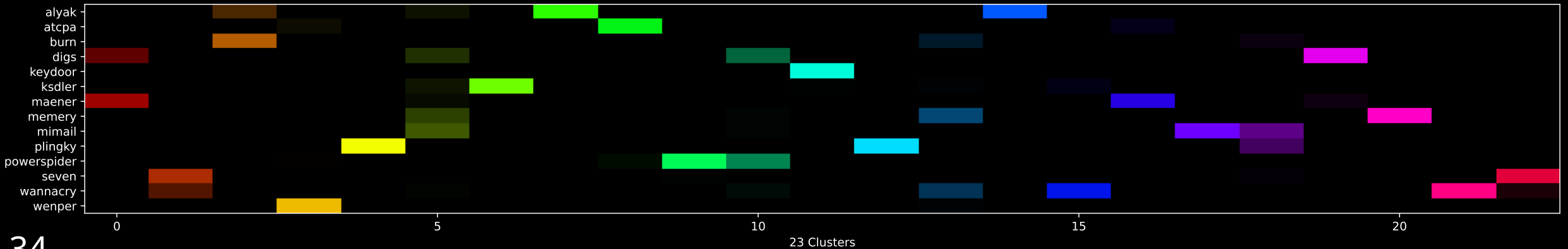- We tried many clustering algorithms, and finally chose spectral clustering

# RESULTS?



TSNE(2) projection of 1321 samples with 50 Metagraphs features built from metalib optimized by the 'Term Frequency - Inverse Sample Frequency' metric with TF-ISF representation

TSNE(2) projection of Spectral Clustering of 1321 samples for 23 clusters with 50 Metagraphs features on metalib optimized by the 'Term Frequency - Inverse Sample Frequency' metric with TF-ISF representation

Ground truth:
- alyak (100)
- atcpa (98)
- burn (89)
- digs (99)
- keydoor (99)
- ksdler (99)
- maener (99)
- memery (100)
- mimail (97)
- plingky (46)
- powerspider (98)
- seven (99)
- wannacry (99)
- wenper (99)

Predictions for K = 23 clusters
RI = 94.24%
ARI = 56.09%

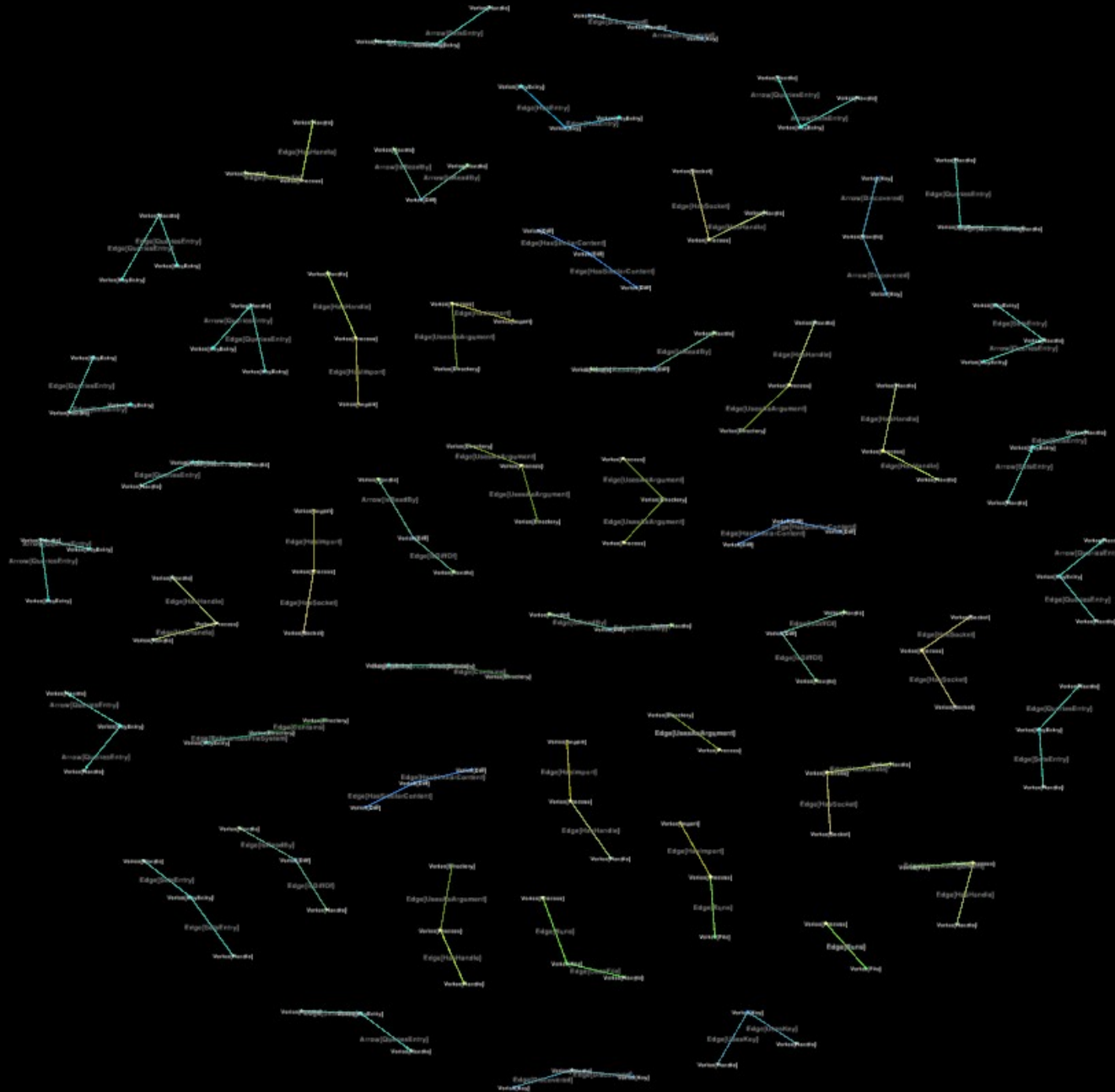Correlation between 14 labels and 23 clusters

23 Clusters

# GOALS

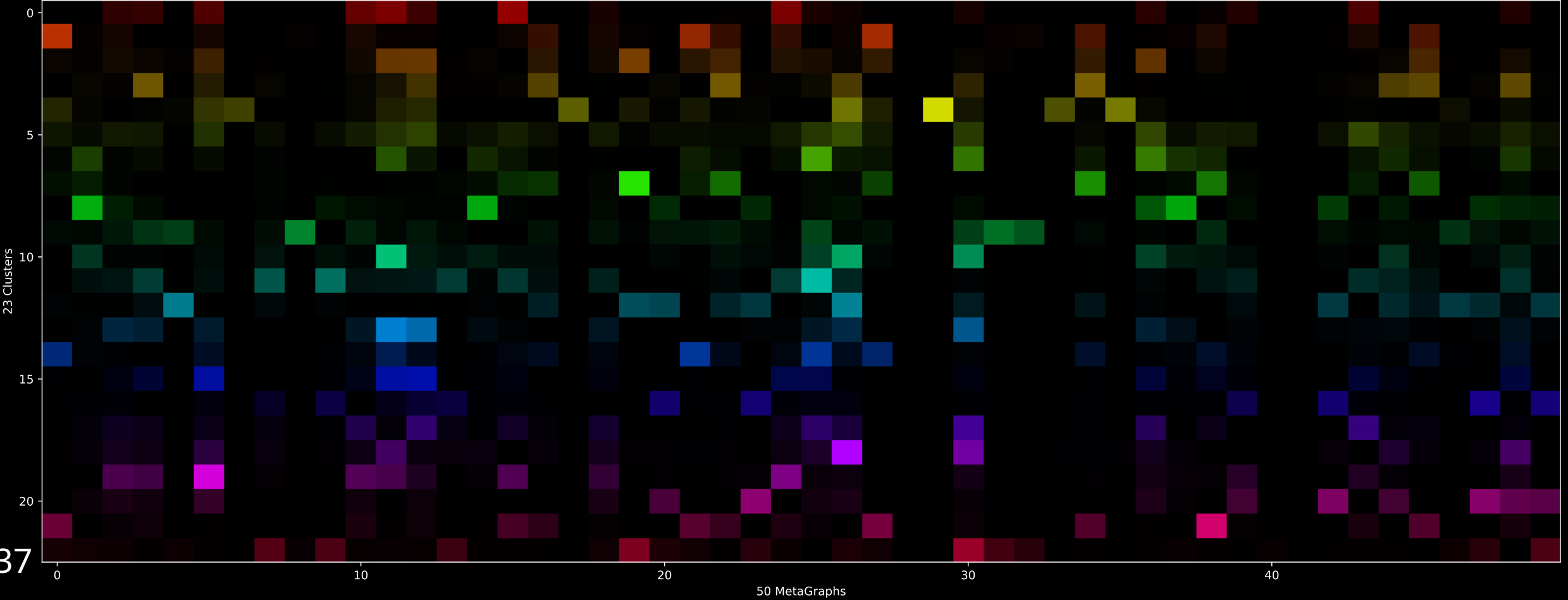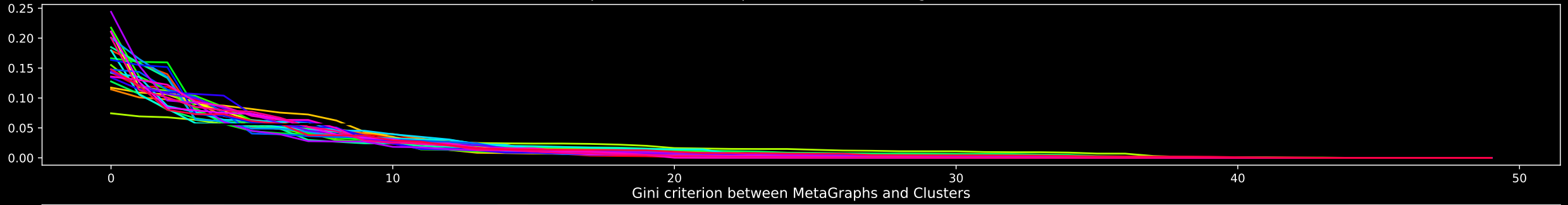Given a unlabelled dataset of malware samples, how can we:

- Recover clusters of behaviors that hopefully match the families/classes → encouraging results

- Recover the behavioral patterns characteristic of each of these clusters
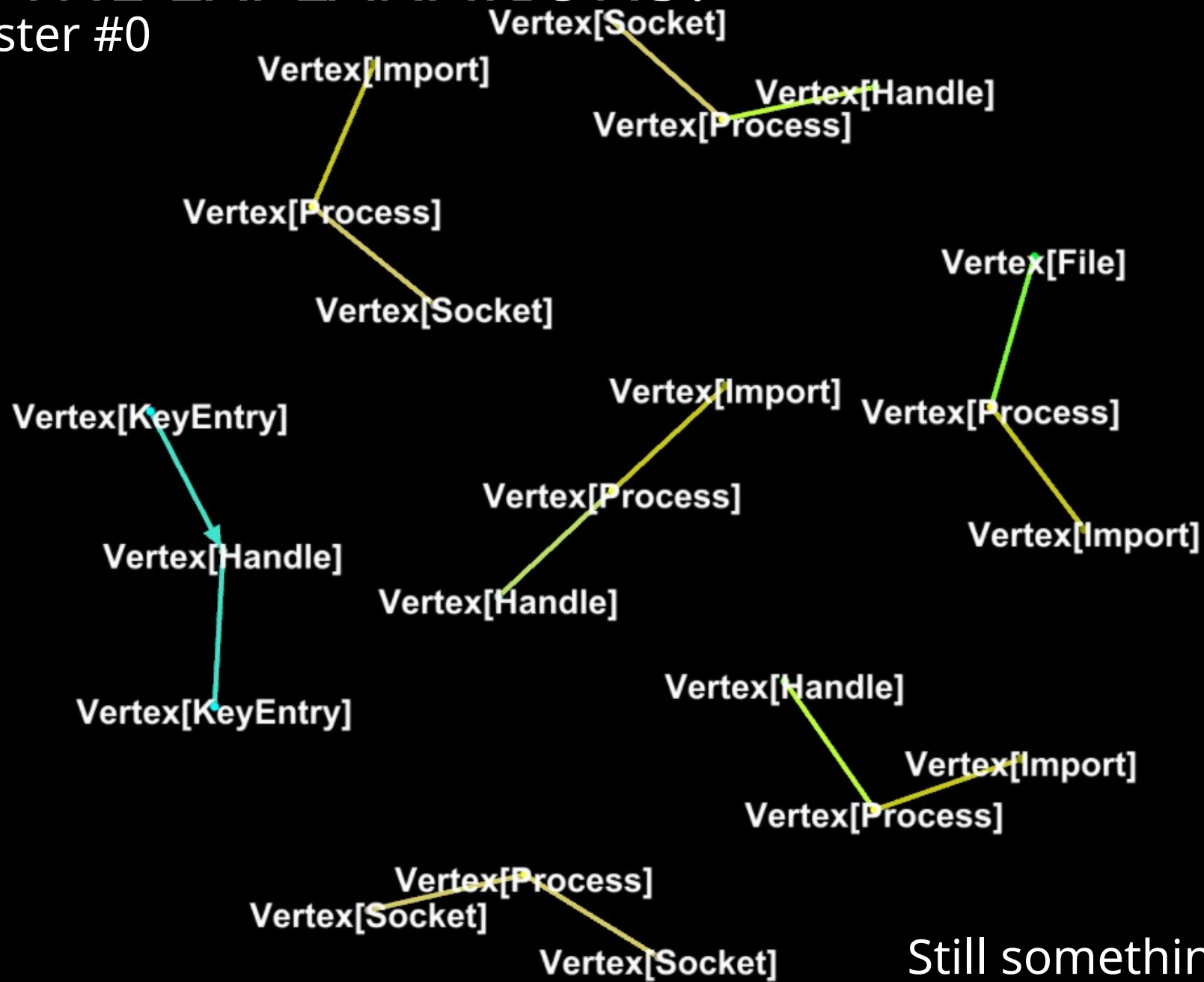
# AND THE EXPLANATIONS?

MetaGraph feature relevance per cluster. Selection weight threshold is 5.6%.



Gini criterion between MetaGraphs and Clusters

# AND THE EXPLANATIONS?
## For cluster #0

Still something to improve...

# NEXT STEPS

- This is still a work in progress
- Proper statistical analysis
    - More iterations to learn more complex metagraphs (with random selection)
    - Clustering algorithm comparison
    - Average accuracy metrics estimation
- Comparison to SOTA of heterogeneous graph pattern mining
- Better metagraph generation rules with better explainability

# CONCLUSION

- BAGUETTE encapsulates the behavioral information of malware samples in a reasonable scale.

- Metagraphs allows to make advanced searches through a dataset of BAGUETTES.

- BAGUETTES give a visualization advantage when analyzing malware samples.

- From BAGUETTE, CROISSANT can:
  - Learn explainable behavioral signatures
  - Can differentiate malware families
  - No need of labelled data