# Attacking the malware with AI

## Where the finest concepts of Data Science & Cybersecurity meet

# Introduction

- Master's in Cybersecurity @ *Georgia Institute of Technology* (USA)

- Security Engineer @ *Trade Republic* (Berlin)

- Cloud Security & DevSecOps

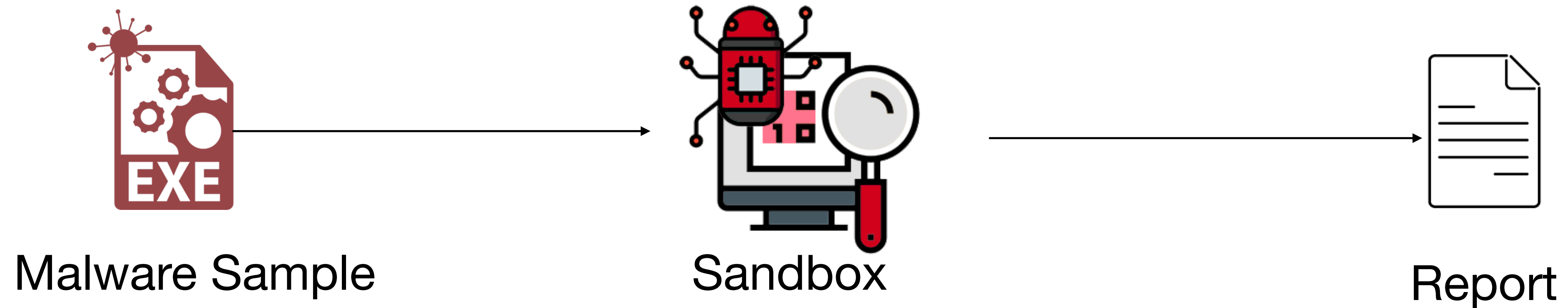- Artificial Intelligence & Privacy

# Agenda

1. Malware Analysis: Sandboxing - How does a Sandbox work?
2. Elements of Machine Learning (Classification & Clustering)
3. Malheur Framework - Explained step by step
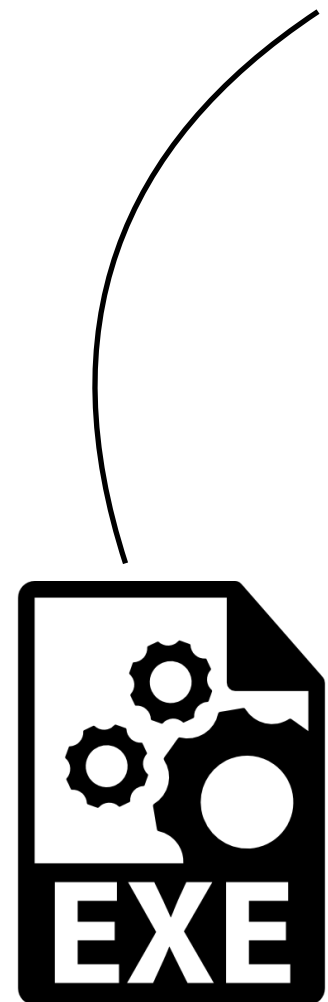4. Q/A

# What is a sandbox?
## And how does it work?

- A sandbox is an **isolated environment** in which malware can be safely executed, in order to **study & monitor** its behaviour

Malware Sample      Sandbox      Report

# What is classification?
## Supervised Learning

**Classification** is the task of identifying the **category** (a.k.a the class) on which an **observation** falls into
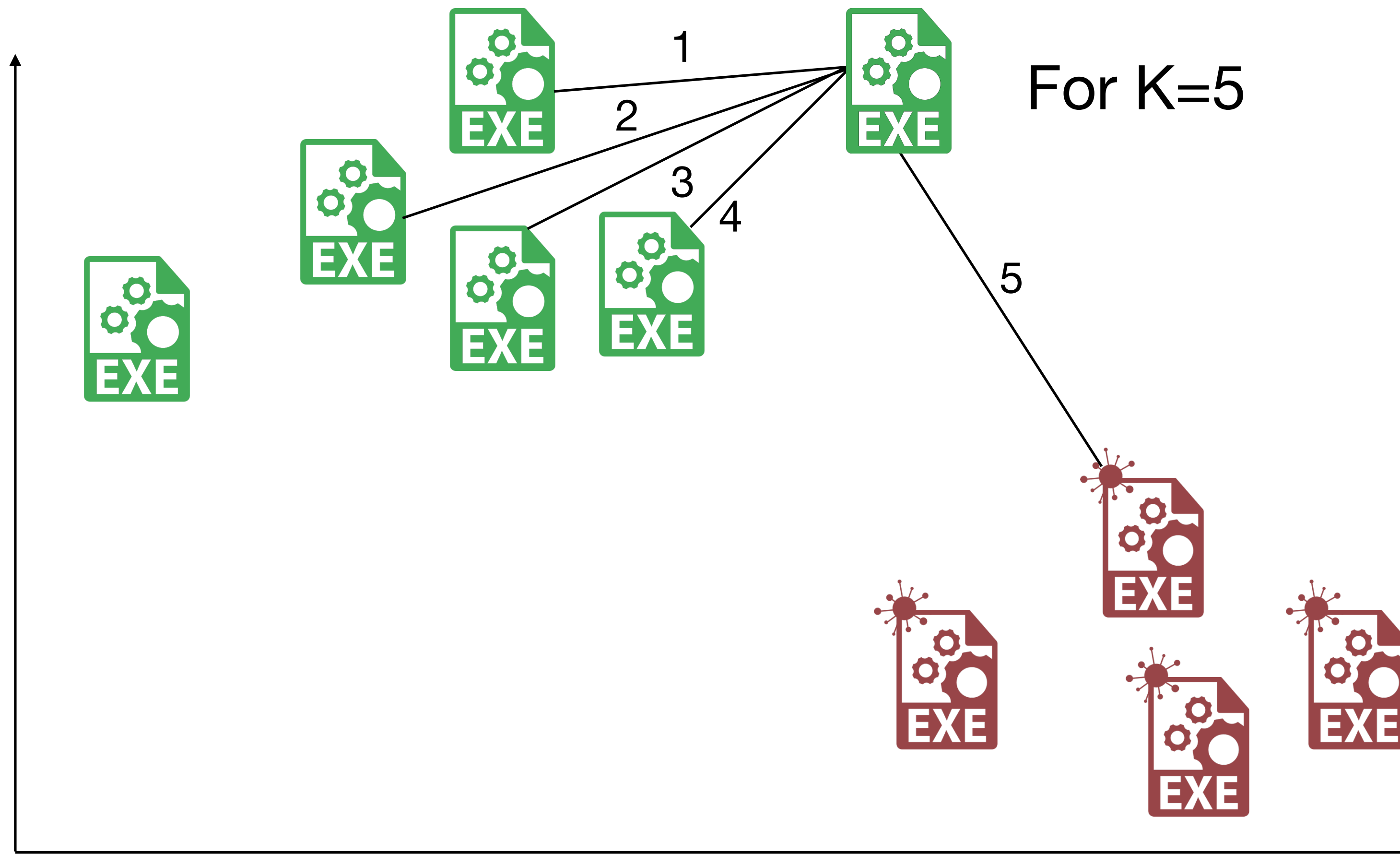
Benign

Malware

# What is classification?
## Supervised Learning



For K=5

There's a folk saying…

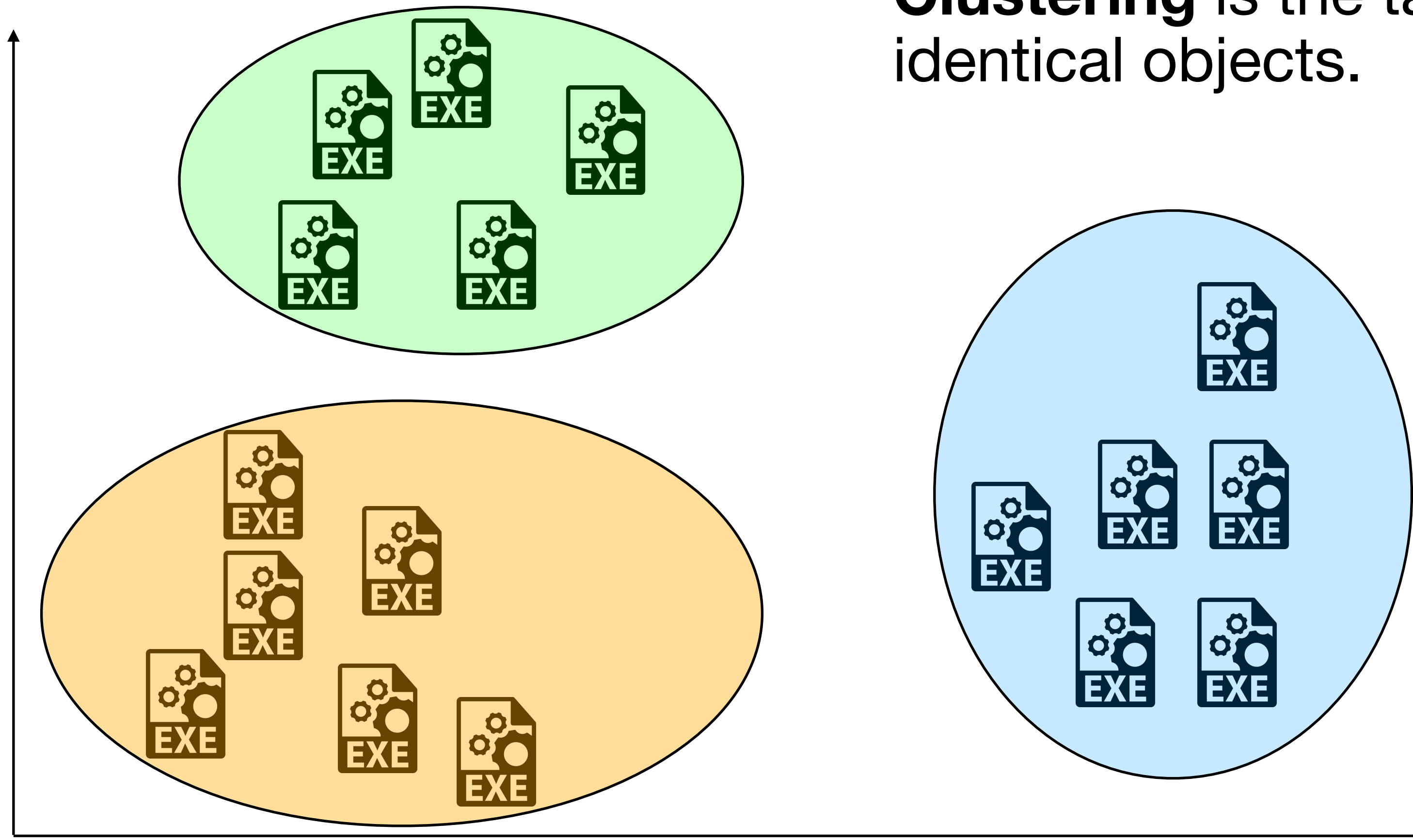"Show me your friends and I will tell you who you are"

There is also an algorithm based on this folk wisdom:

**K-Nearest** ~~Friends~~ **Neighbours**

# What is clustering?
## Unsupervised learning

**Clustering** is the task of forming groups of identical objects.

# Malheur Framework (Rieck et al., 2011)

*"More than 450,000 new malicious software and potentially unwanted applications (PUA) are registered every day"*

*AV-Test Institute, 2022*

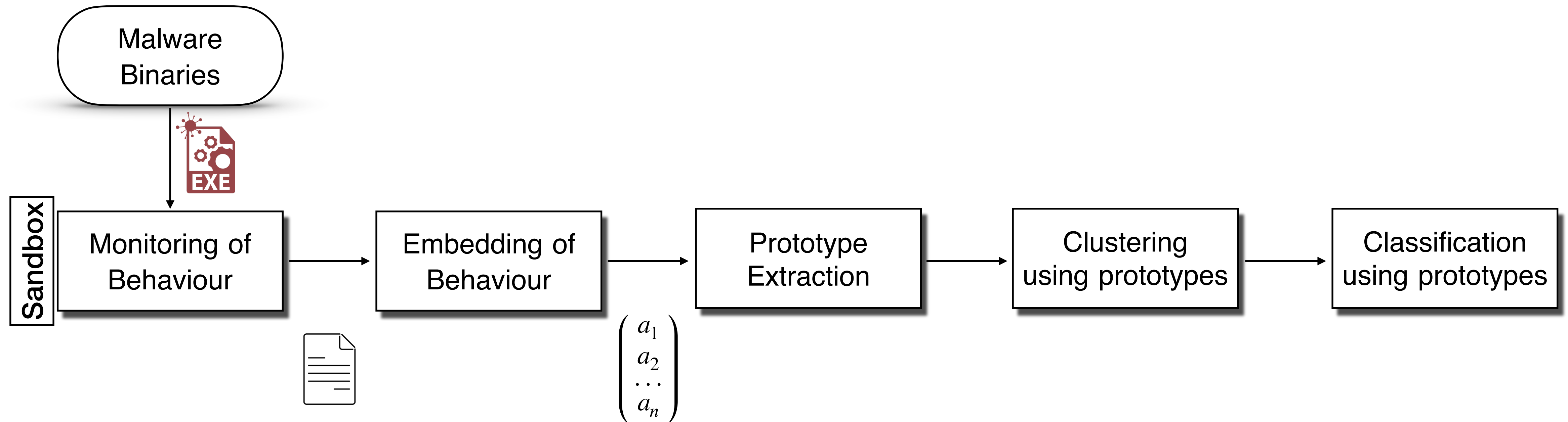*Malheur* is an open-source framework for Malware Analysis with Machine Learning.

It is used to:

1. Automatically discover novel classes of malware

2. Classify unknown malware to known classes

# Malheur Framework Simplified
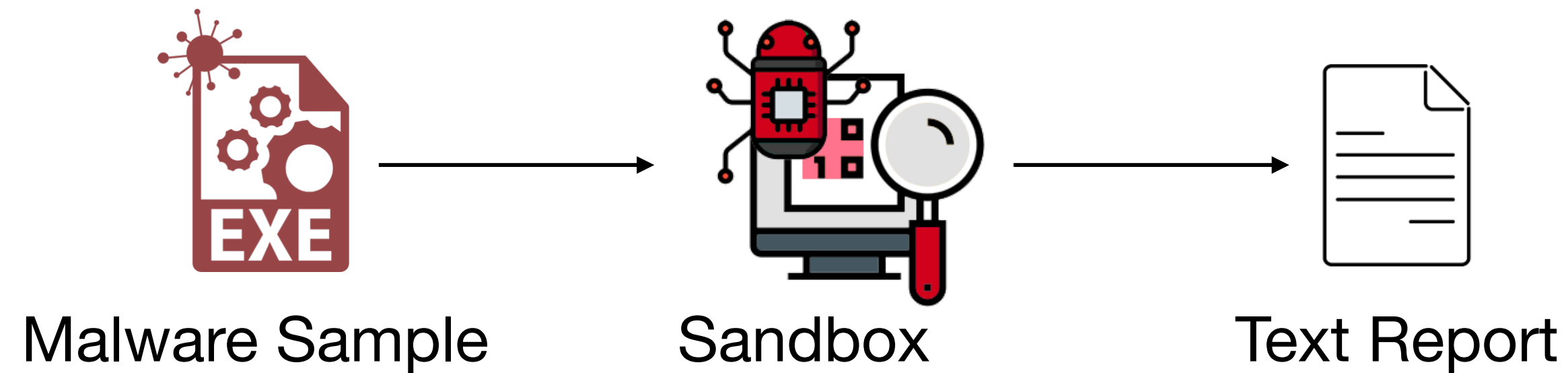## Overview

Malware Binaries

Sandbox

Monitoring of Behaviour → Embedding of Behaviour → Prototype Extraction → Clustering using prototypes → Classification using prototypes

$$\begin{pmatrix} a_1 \\ a_2 \\ \dots \\ a_n \end{pmatrix}$$

# Step 1: Monitoring of the behaviour

Malware Sample → Sandbox → Text Report

```
<load_dll filename="C:\WINDOWS\system32\kernel32.dll"
<get_system_time apifunction="GetLocalTime"/>
...

...

<copy_file filetype="file" srcfile="c:\\malware.exe" dstfile="C:
\system32\csrss.exe"/>
<set_value key="HKEY_LOCAL_MACHINE\CurrentVersion\Run" value="UpDaTer" data="C:
\system32\csrss.exe"/>
...

...

<check_for_debugger apifunction="IsDebuggerPresent"/>
<load_dll filename="C:\WINDOWS\system32\UxTheme.dll"/>
...

...
```
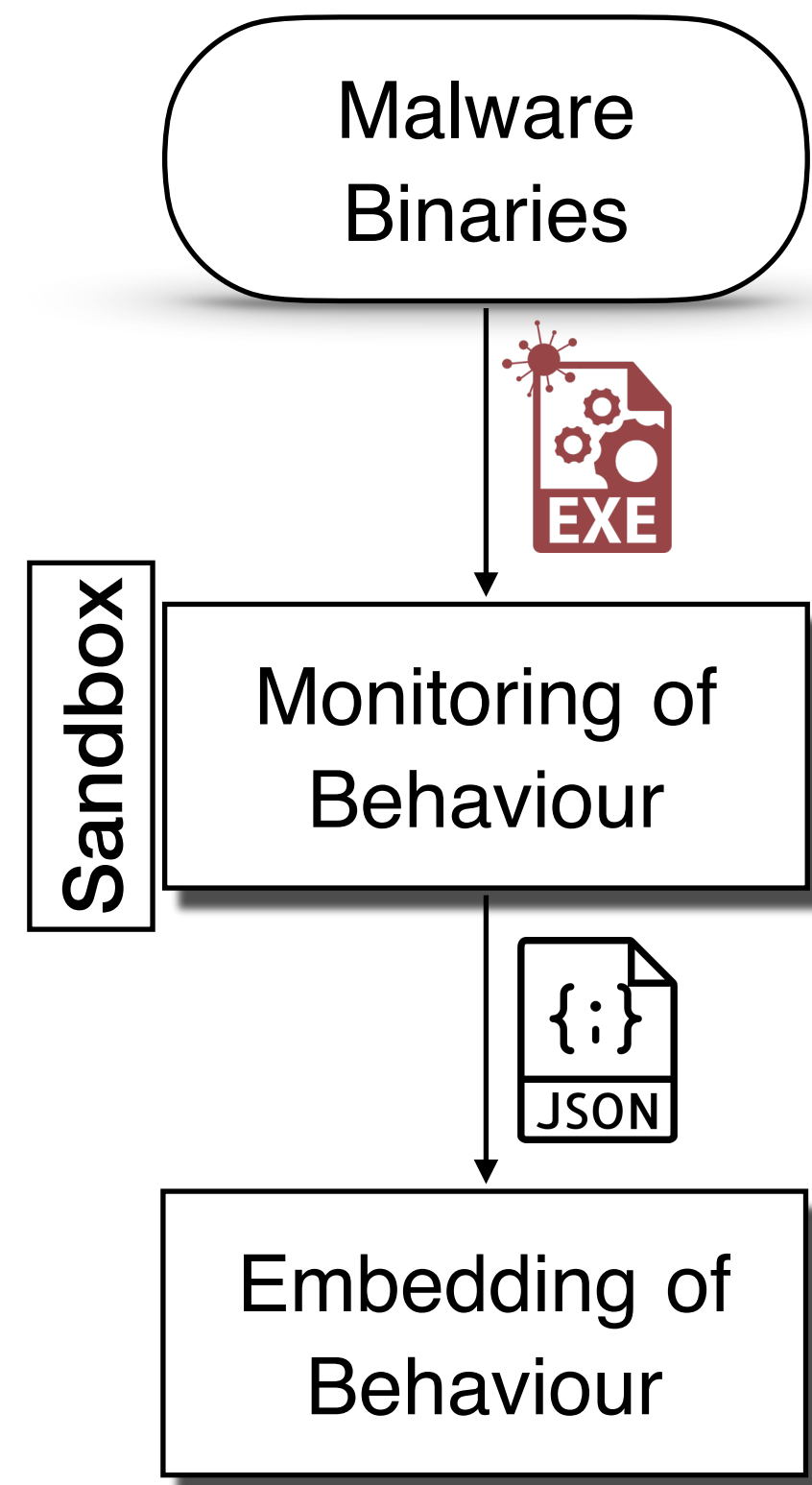
# Step 2: Embedding of the behaviour
## Instruction Q-grams

- Behaviour is often manifested as a **sequence of instructions,** which are formally known as **q-grams**

- For example, a 2-gram:

```
1. <copy_file filetype="file" srcfile="c:\\malware.exe" dstfile="C:
\system32\csrss.exe"/>
2. <set_value key="HKEY_LOCAL_MACHINE\CurrentVersion\Run" value="UpDaTer" data="C:
\system32\csrss.exe"/>
```

```
1. <check_for_debugger apifunction="IsDebuggerPresent"/>
2. <load_dll filename="C:\WINDOWS\system32\UxTheme.dll"/>
```

# Step 2: Embedding of the behaviour

Malware Binaries

Sandbox

Monitoring of Behaviour

JSON

Embedding of Behaviour

The embedding of reports in vector spaces enables us to express the similarity of behaviour *geometrically*

We model the text files into mathematical vectors using an ***embedding function***
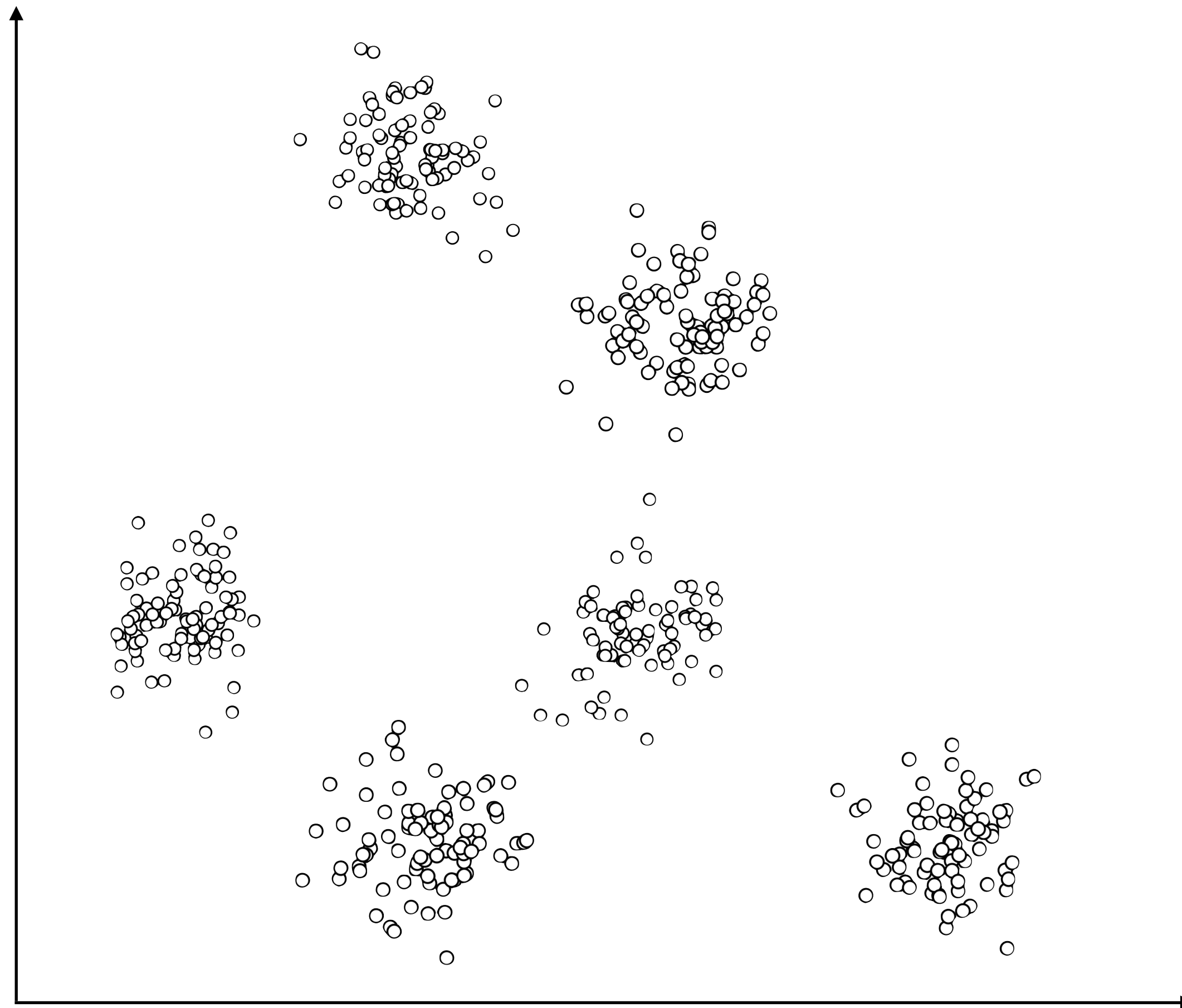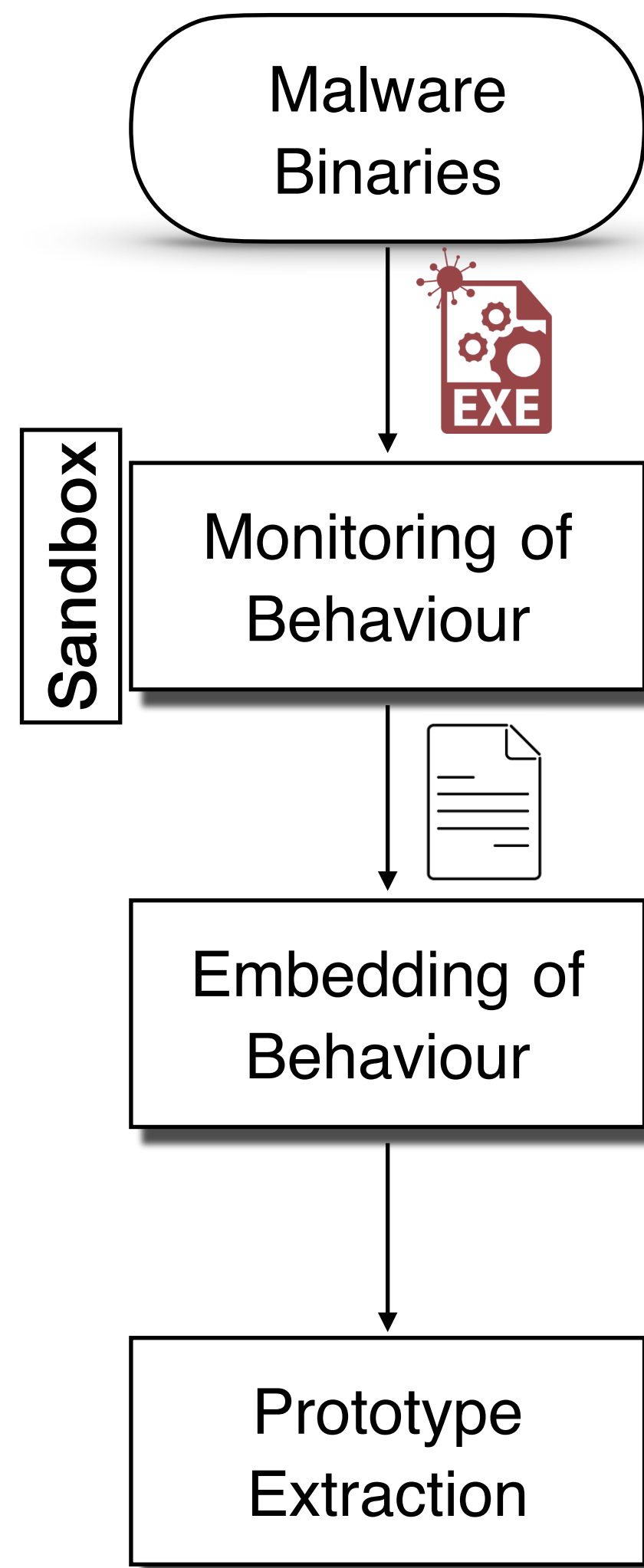
# Step 2: Embedding of the behaviour

Malware
Binaries

Sandbox

Monitoring of
Behaviour

Embedding of
Behaviour

$$\mathcal{S} = \{(a_1, \ldots, a_q) \mid a_i \in \mathcal{A} \text{ with } 1 \leq i \leq q\},$$

$$\phi(x) = (\phi_s(x))_{s \in S} \quad \text{with} \quad \phi_s(x) = \begin{cases} 1 & \text{if report x contains q-grams S,} \\ 0 & otherwise \end{cases}$$

$$\phi({}'A_1 \quad A_2 \quad A_1 \quad A_2') \rightarrow \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix} \quad \begin{array}{cc} {}'A_1 & A_1' \\ {}'A_1 & A_2' \\ {}'A_2 & A_1' \\ {}'A_2 & A_2' \end{array}$$
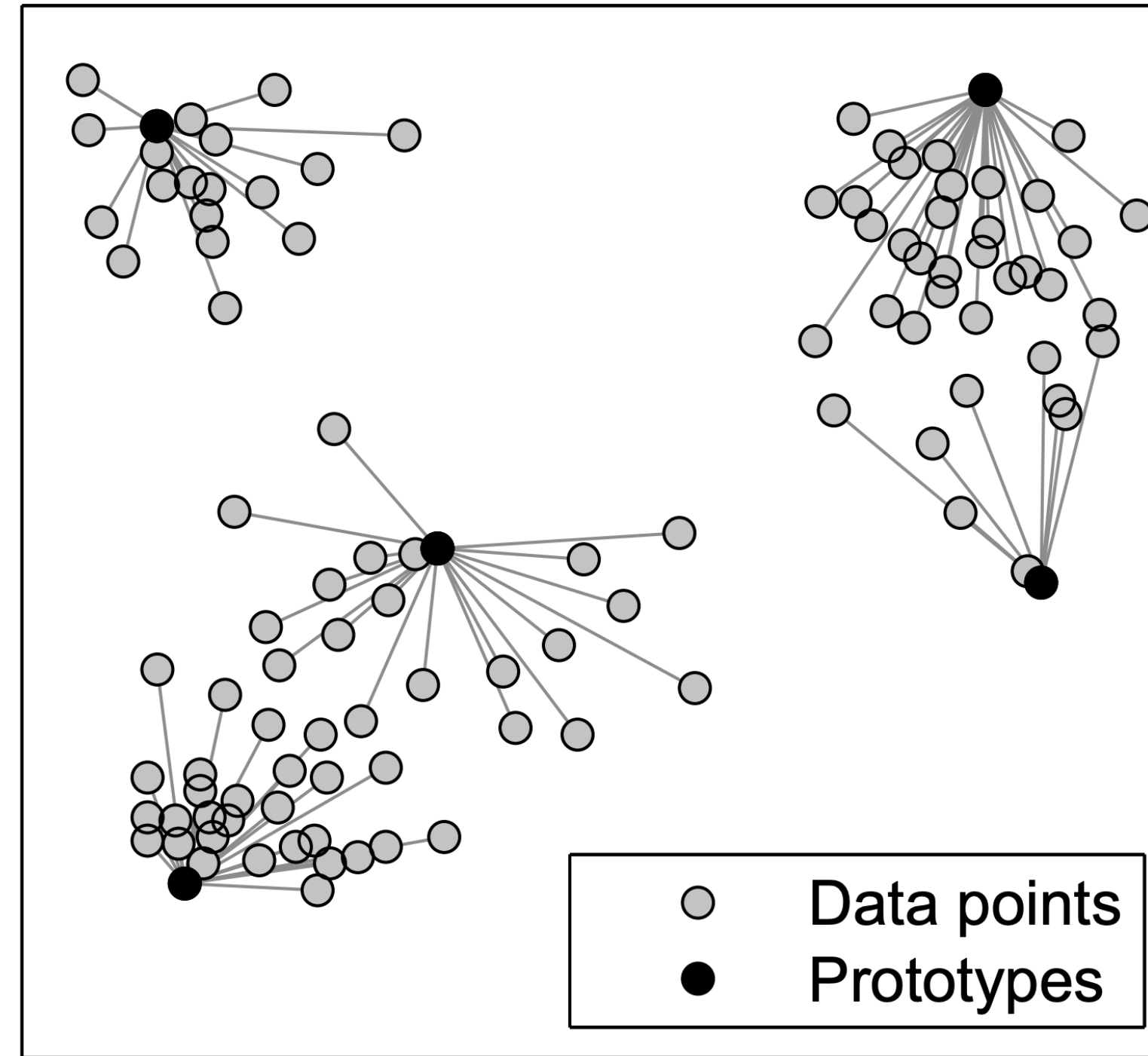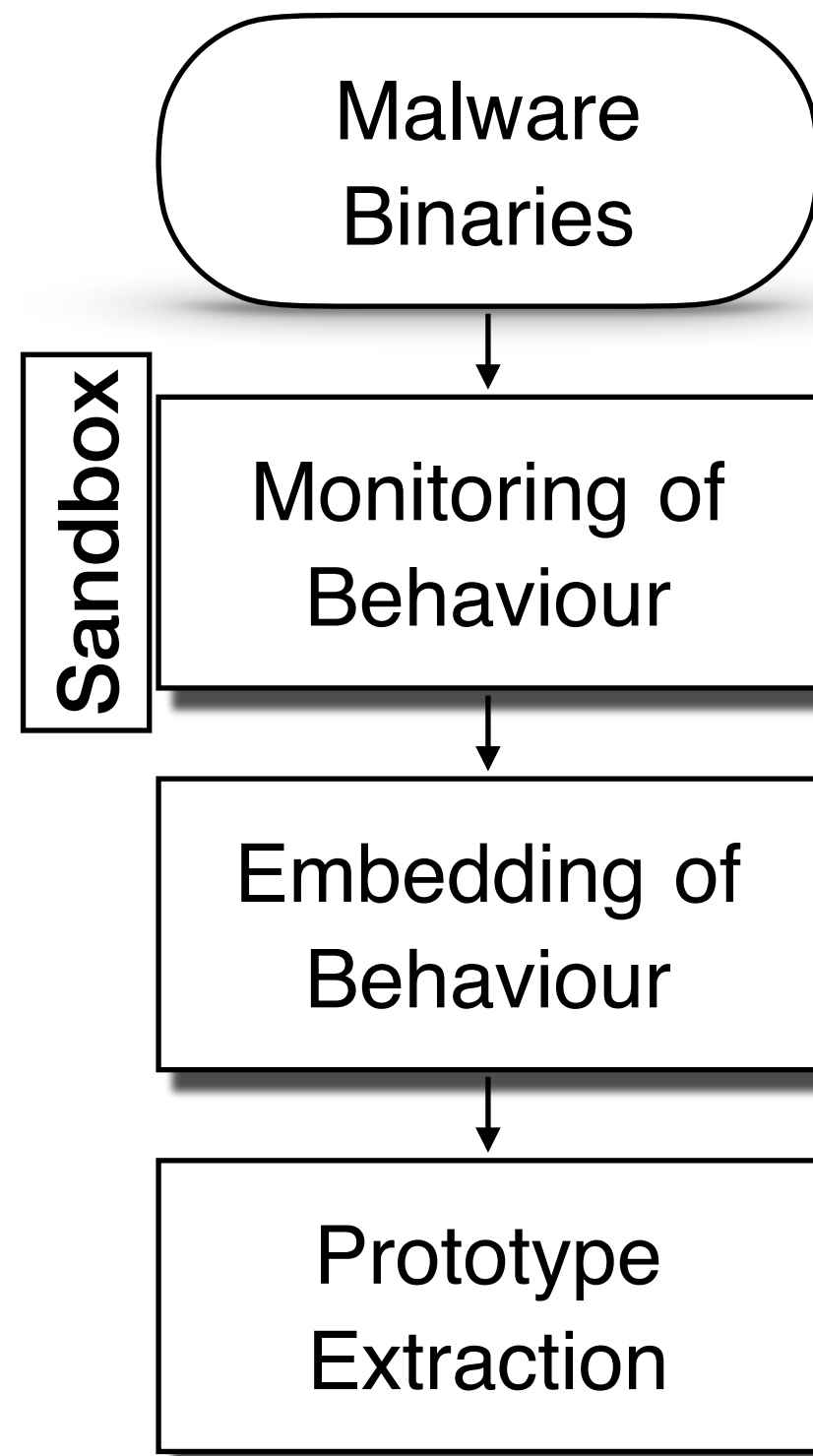
# Step 2: Embedding of Behaviour

Malware
Binaries

Sandbox

Monitoring of
Behaviour

Embedding of
Behaviour

Prototype
Extraction

- Identical reports *"form dense clouds in the vector space"*

# Computational Complexity

# Step 3: Prototype Extraction

Malware Binaries

Sandbox

Monitoring of Behaviour

Embedding of Behaviour

Prototype Extraction

Data points
Prototypes

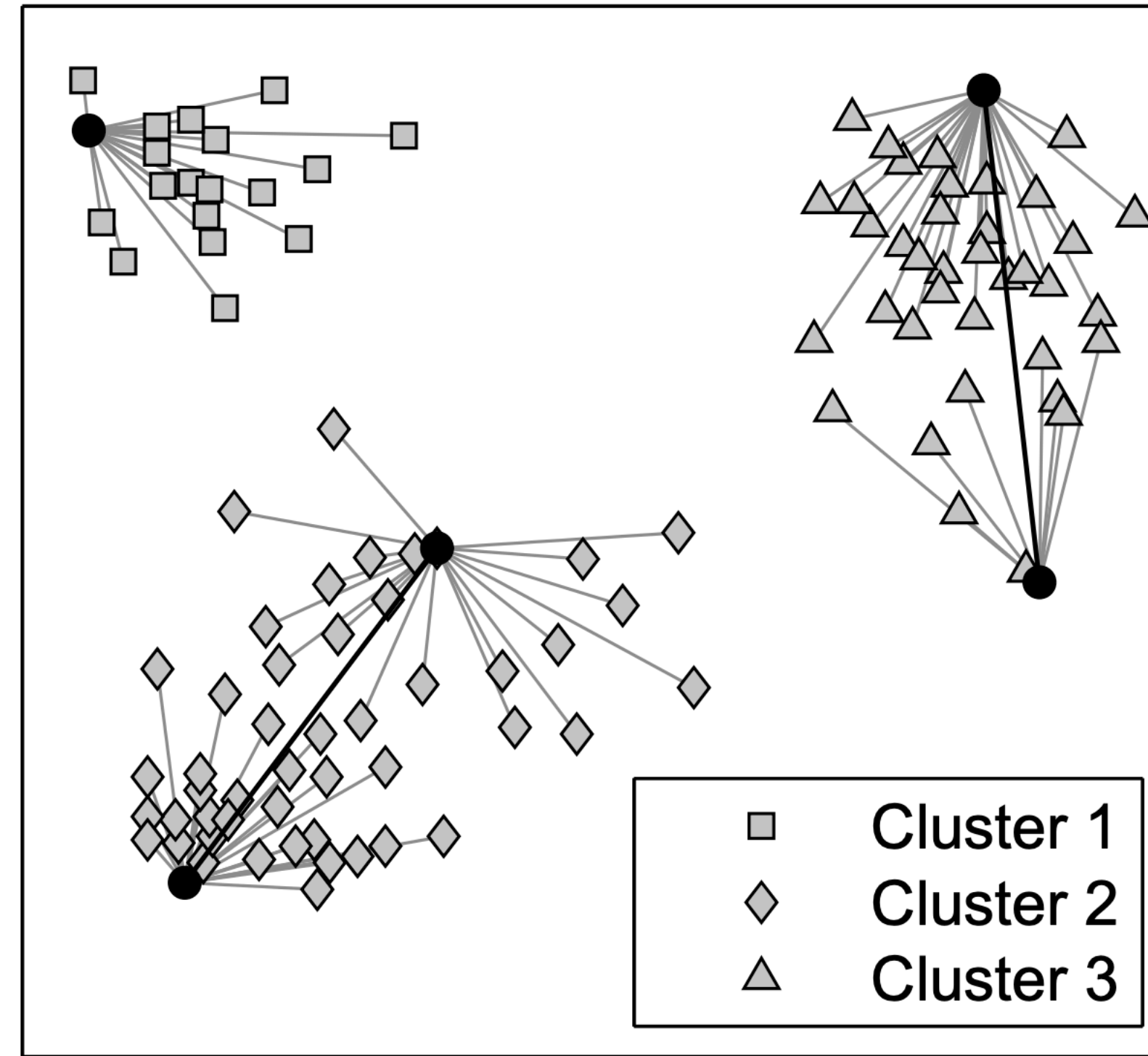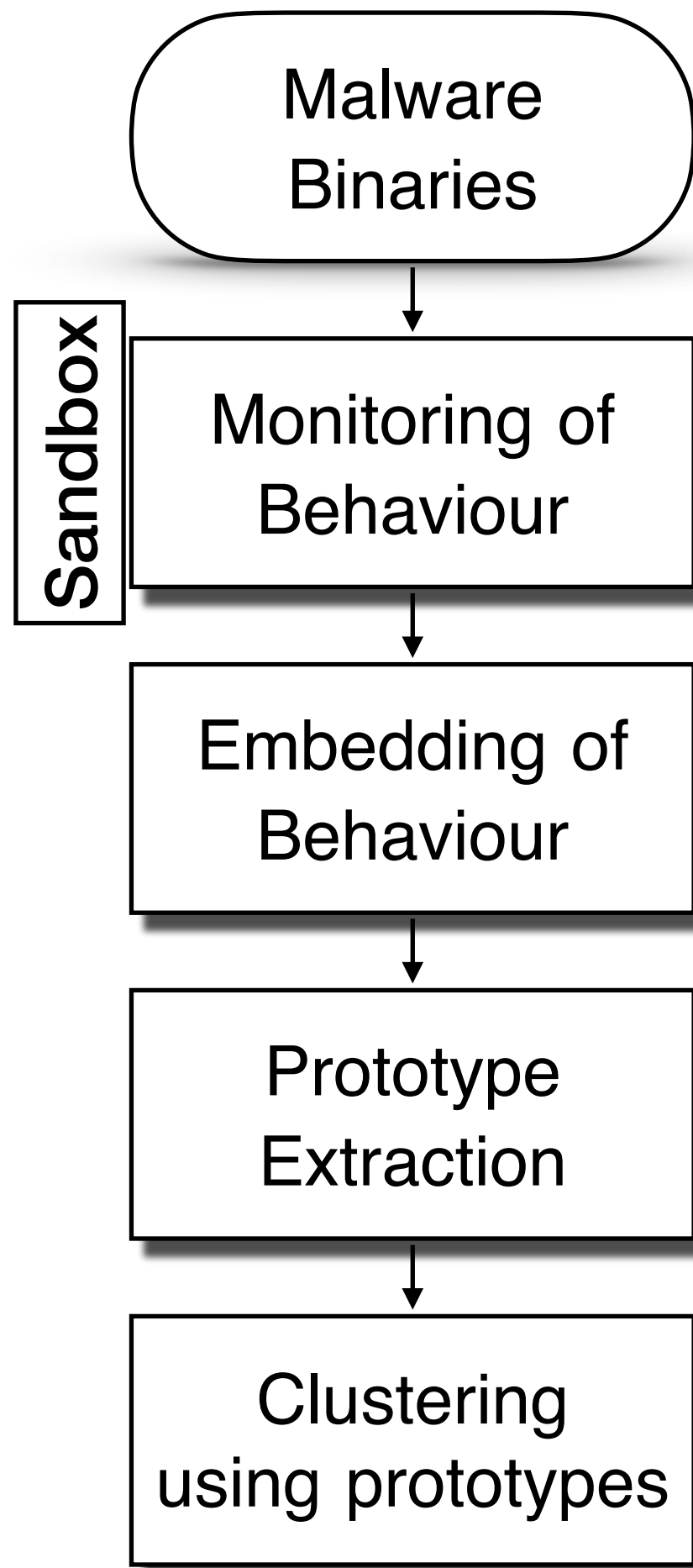**Algorithm 1** Prototype extraction

$O(k \cdot n)$

1: $prototypes \leftarrow \varnothing$
2: $distance[x] \leftarrow \infty$ for all $x \in reports$
3: **while** $\max(distance) > d_p$ **do**
4:     choose $z$ such that $distance[z] = \max(distance)$
5:     **for** $x \in reports$ and $x \neq z$ **do**
6:         **if** $distance[x] > ||\hat{\varphi}(x) - \hat{\varphi}(z)||$ **then**
7:             $distance[x] \leftarrow ||\hat{\varphi}(x) - \hat{\varphi}(z)||$
8:     add $z$ to $prototypes$

`k: # of prototypes`
`n: # of reports`

(Rieck et al., 2011)

# Step 4: Clustering using Prototypes

Malware Binaries

Sandbox

Monitoring of Behaviour

Embedding of Behaviour

Prototype Extraction

Clustering using prototypes



□ Cluster 1
◇ Cluster 2
△ Cluster 3

"Once a clustering has been determined on the prototypes, it is propagated to the original reports"
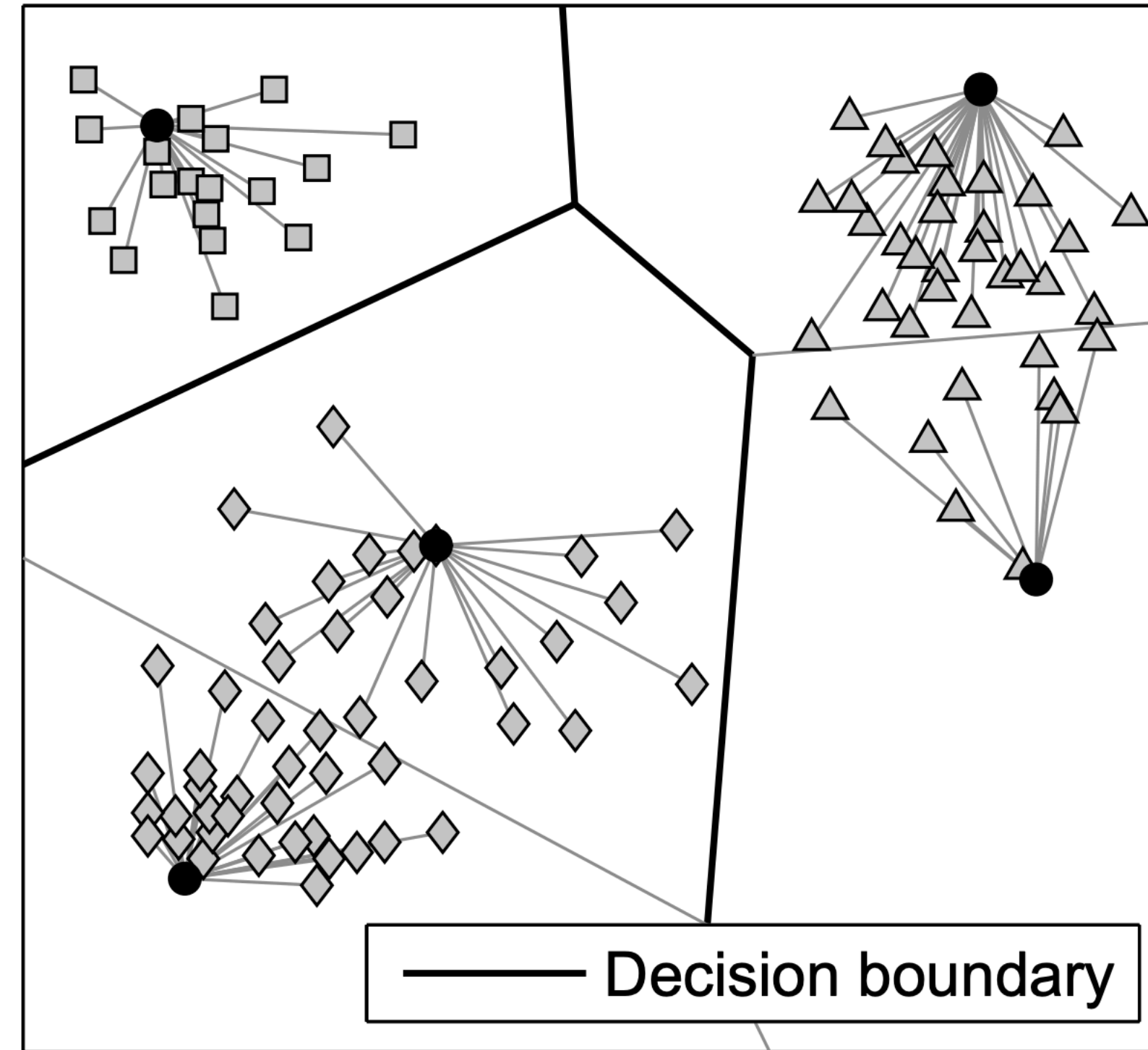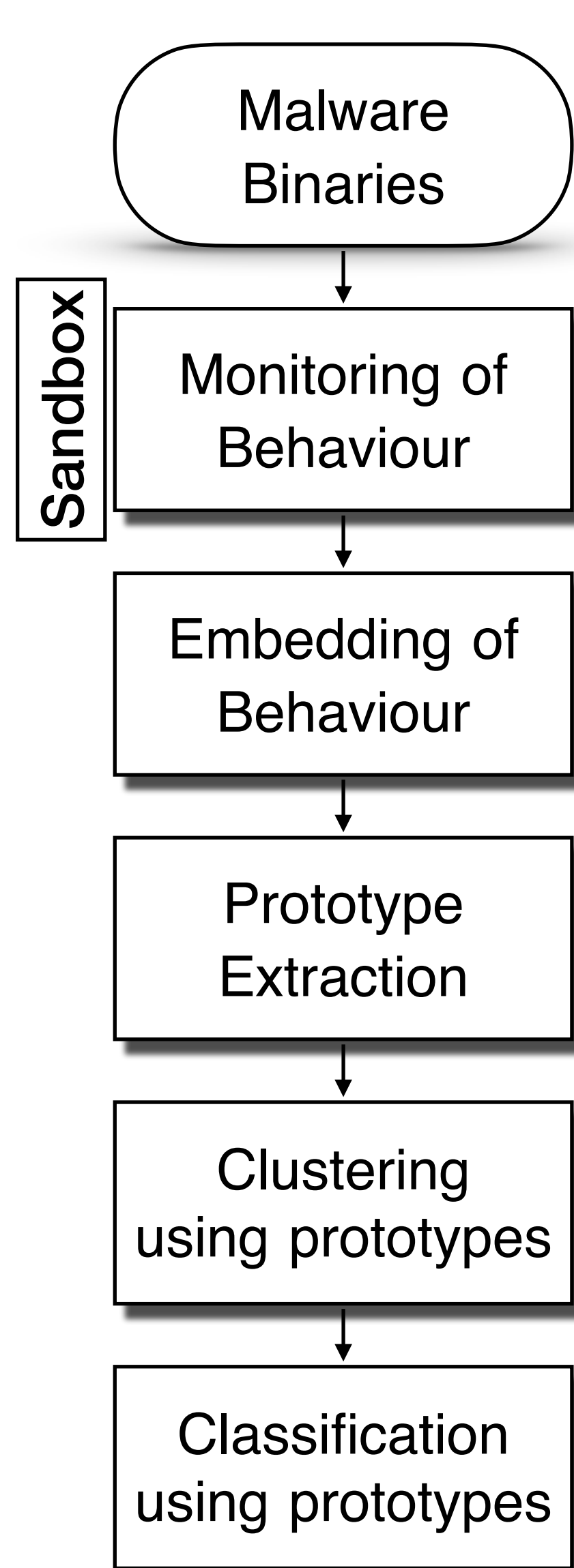
**Algorithm 2** Clustering using prototypes

$O(k^2 \cdot log(k) + n)$

1: **for** $z, z' \in$ *prototypes* **do**
2:     $distance[z, z'] \leftarrow ||\hat{\varphi}(z) - \hat{\varphi}(z')||$
3: **while** $\min(distance) < d_c$ **do**
4:     merge clusters $z, z'$ with minimum $distance[z, z']$
5:     update *distance* using complete linkage
6: **for** $x \in$ *reports* **do**
7:     $z \leftarrow$ nearest prototype to $x$
8:     assign $x$ to cluster of $z$
9: reject clusters with less than $m$ members

k: # of prototypes
n: # of reports

"clusters with fewer than **m** members are rejected and kept for later incremental analysis"

(Rieck et al., 2011)

# Step 5: Classification using Prototypes



```
Malware
Binaries
```

Sandbox

```
Monitoring of
Behaviour
```

```
Embedding of
Behaviour
```

```
Prototype
Extraction
```

```
Clustering
using prototypes
```

```
Classification
using prototypes
```

— Decision boundary

**Algorithm 3** Classification using prototypes          $O(k \cdot n)$
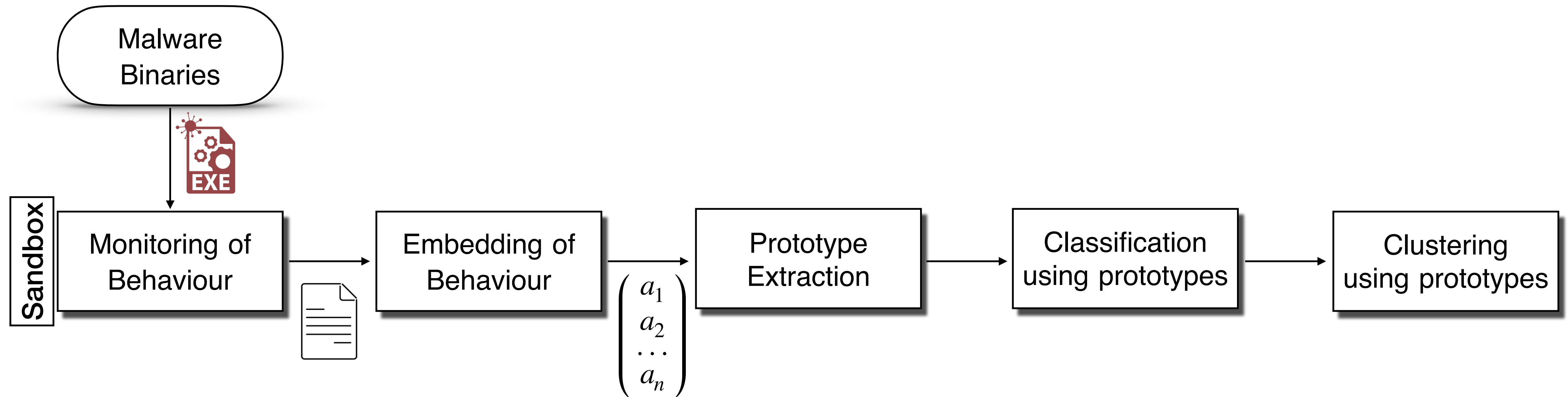
1: **for** $x \in$ *reports* **do**
2:     $z \leftarrow$ nearest prototype to $x$
3:     **if** $||\hat{\varphi}(z) - \hat{\varphi}(x)|| > d_r$ **then**
4:         reject $x$ as unknown class
5:     **else**
6:         assign $x$ to cluster of $z$

k: # of prototypes
n: # of reports

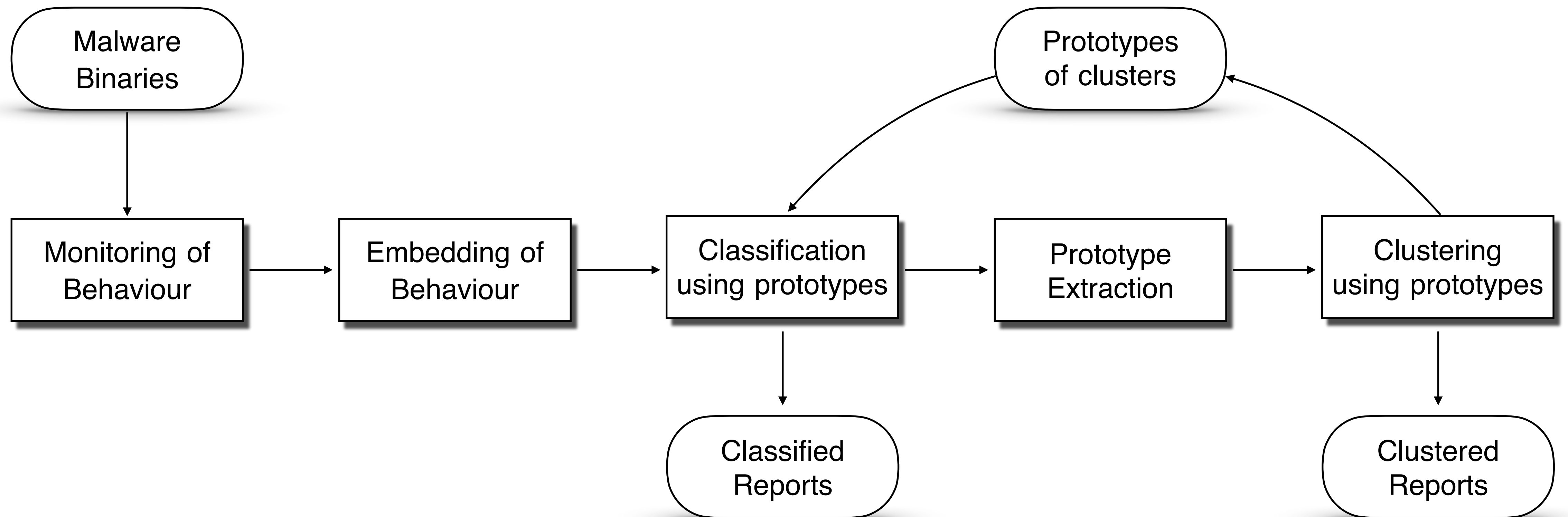"Nearest prototype classification resembles an efficient alternative to costly k-nearest neighbor methods"

(Rieck et al., 2011)

# Malheur Framework Simplified
## Overview

Malware Binaries

Sandbox

Monitoring of Behaviour

Embedding of Behaviour

$$\begin{pmatrix} a_1 \\ a_2 \\ \dots \\ a_n \end{pmatrix}$$

Prototype Extraction

Classification using prototypes

Clustering using prototypes

# Malheur Framework
## Overview

# Incremental Analysis

---

**Algorithm 4** Incremental Analysis

---

1: *rejected* ← ∅, *prototypes* ← ∅
2: **for** *reports* ← data source ∪ *rejected* **do**
3:     classify *reports* to known clusters using *prototypes*      ▷ see Algorithm 3
4:     extract prototypes from remaining *reports*      ▷ see Algorithm 1
5:     cluster remaining *reports* using prototypes      ▷ see Algorithm 2
6:     *prototypes* ← *prototypes* ∪ prototypes of new clusters
7:     *rejected* ← rejected reports from clustering

---

(Rieck et al., 2011)

# Attacking the malware with AI

## Where the finest concepts of Data Science & Cybersecurity meet

- The talk is based on the paper *"Automatic analysis of malware behavior using machine learning"* (Rieck et al., 2011)

Thanks to:

- Sneha Rajguru

- Bsides Munich

# Q/A