

Breaking Into The AI-protected Realm

A Practical Session on (pen)testing Embedded Targets with AI-based
EDR

> Whoarewe

Lalie:



- > PhD Student at CEA-Leti
- > Studies energetic networks' cybersecurity using HW/SW related signals

- > Nix Lover 
- > ENSIMAG alumni

- > Uses AI to filter and identify threats in OS signals trace

Ulysse:

- > PhD Student at CEA-Leti
- > Studies fuzzing with side-channel
- > Has studied the subject of interest in the past
- > Parisian bobo exiled to the mountains
- > Uses statistical tools to perform signal processing

> Who are we

Lalie:

- > PhD Student at CEA-Leti
- > Studies energetic networks' cybersecurity using HW/SW related signals
- > Nix Lover
- > ENSIMAG alumni
- > Use IA to filter and identify threats in OS signals trace

Ulysse:



- > PhD Student at CEA-Leti
- > Studies fuzzing with side-channel
- > Has studied the subject of interest in the past
- > Parisian bobo exiled to the mountains
- > Uses statistical tools to perform signal processing

> What will u do today?

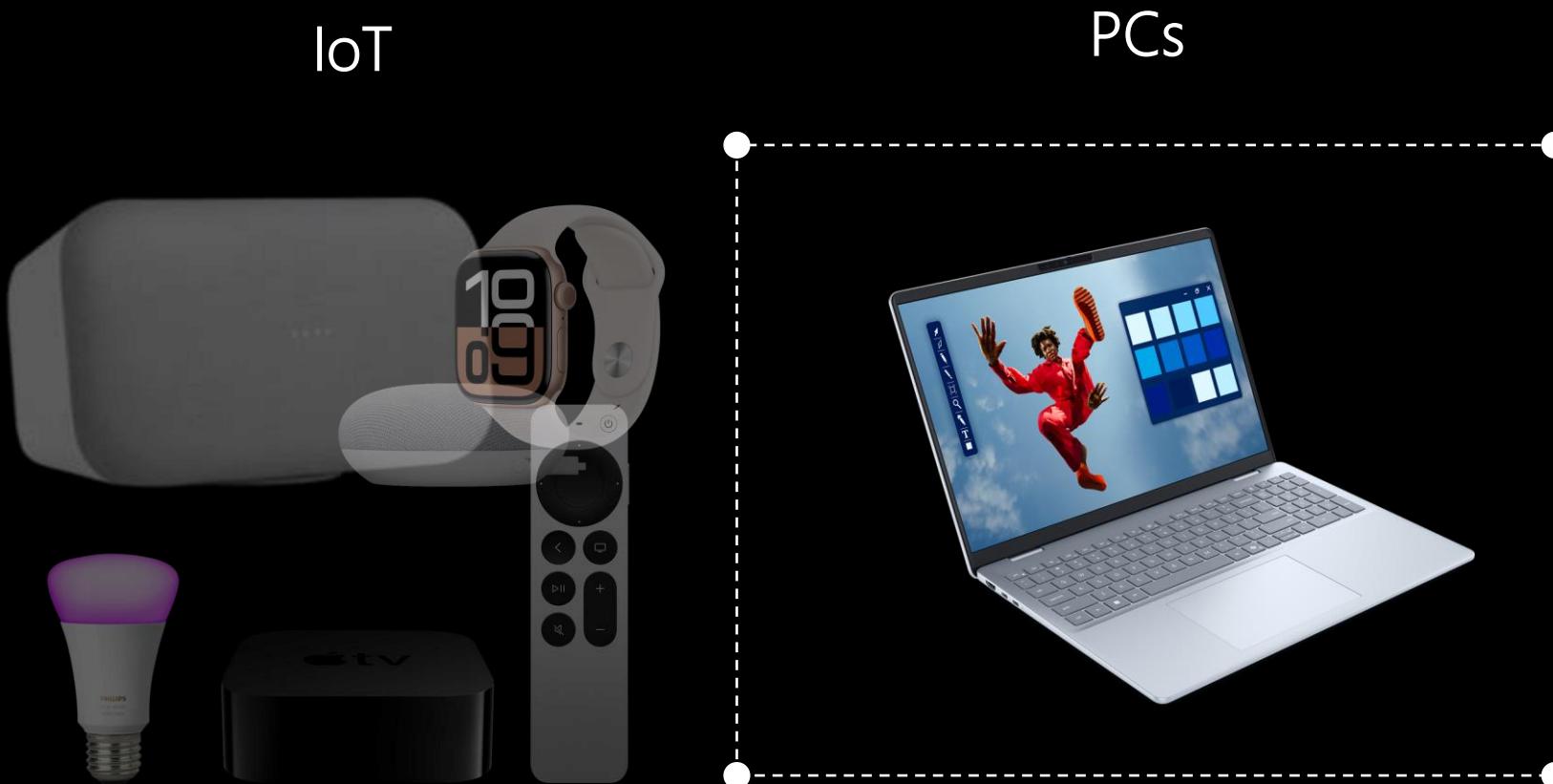
- > Understand why Network-based Detection Systems (NIDS) are not always up to the task
- > Understand what a Host-based Intrusion Detection System (HIDS) is
- > Introduction to some simple statistical tools to perform anomaly detection
- > Try to break into the target without triggering the detection mechanism

> What are IRL current targets?

IoT



> What are IRL current targets?



> What are IRL current targets?

IoT



PCs



Uncrewed vehicles



> What are IRL current targets?

Uncrewed vehicles



> What are IRL current targets?

Uncrewed vehicles



Servers



> What are IRL current targets?

Uncrewed vehicles



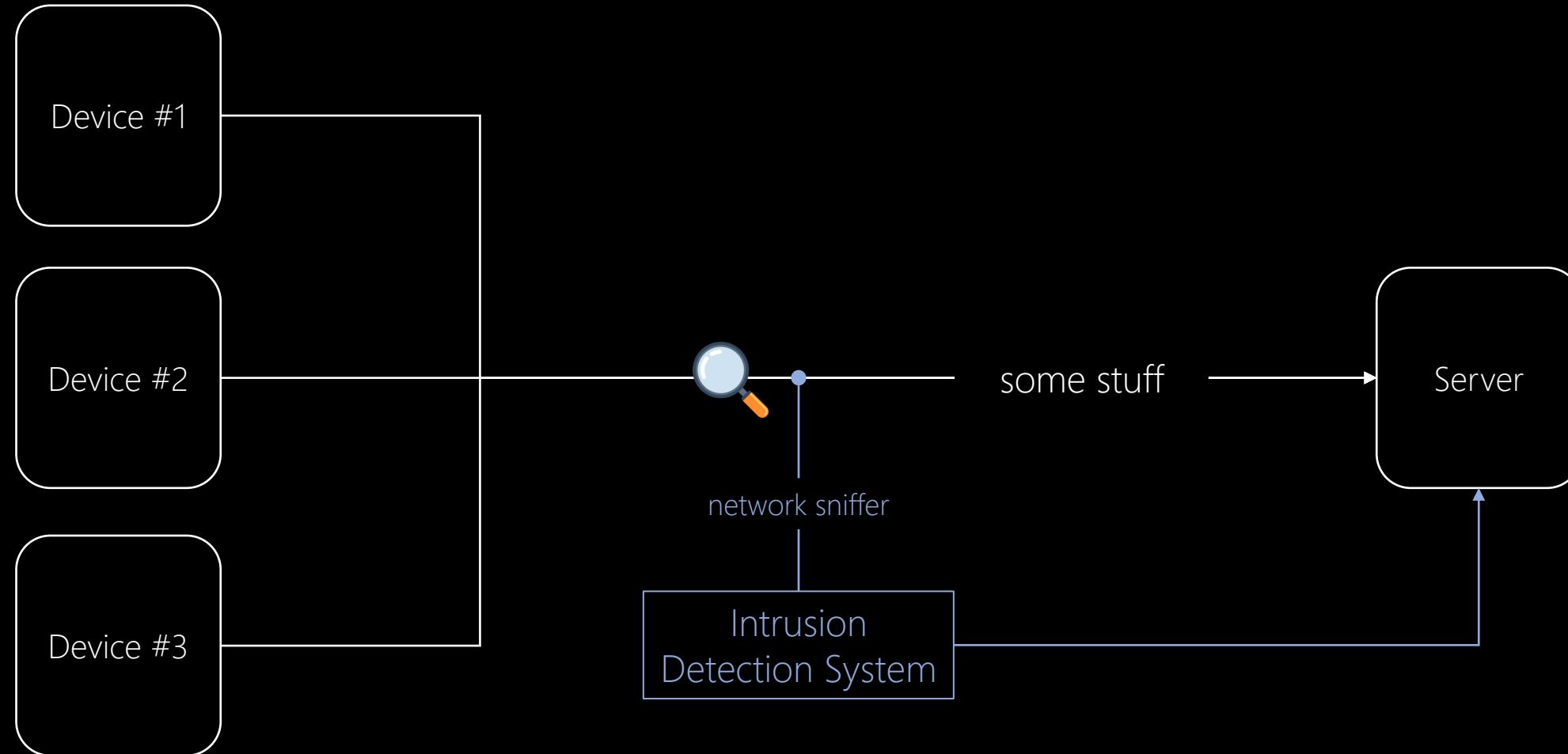
Servers



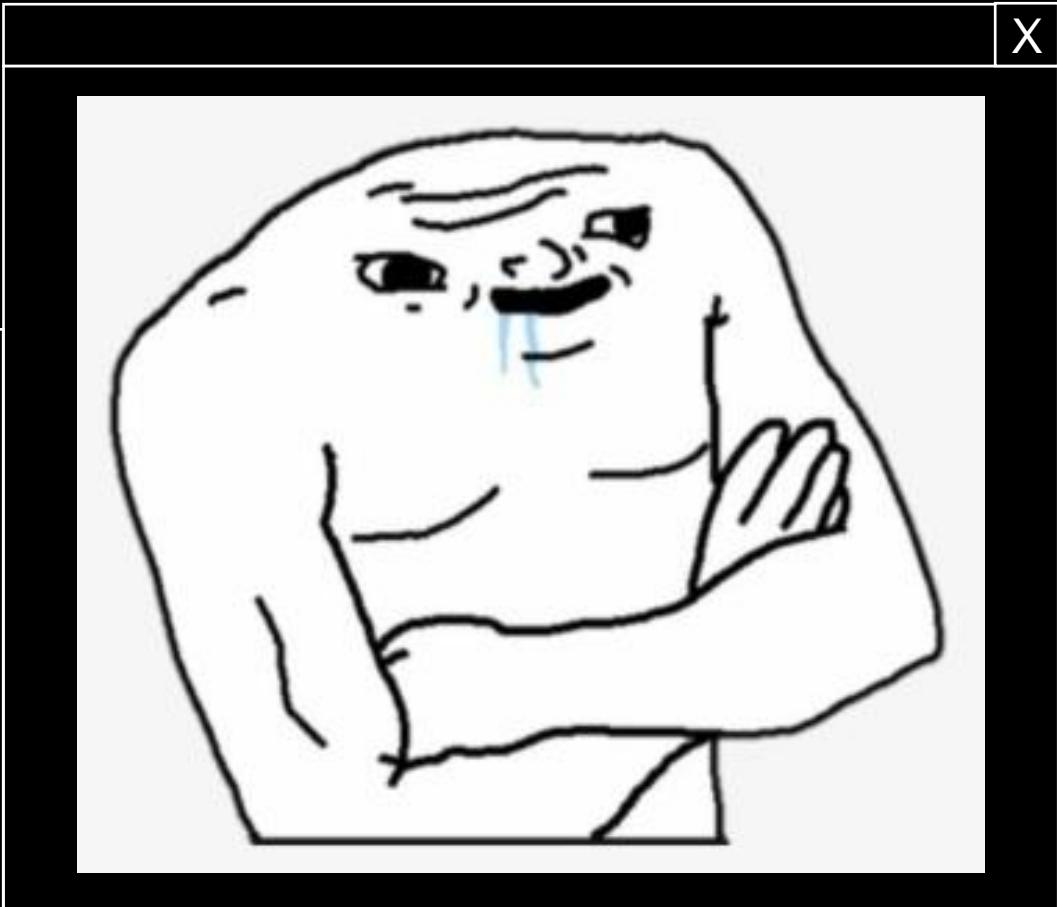
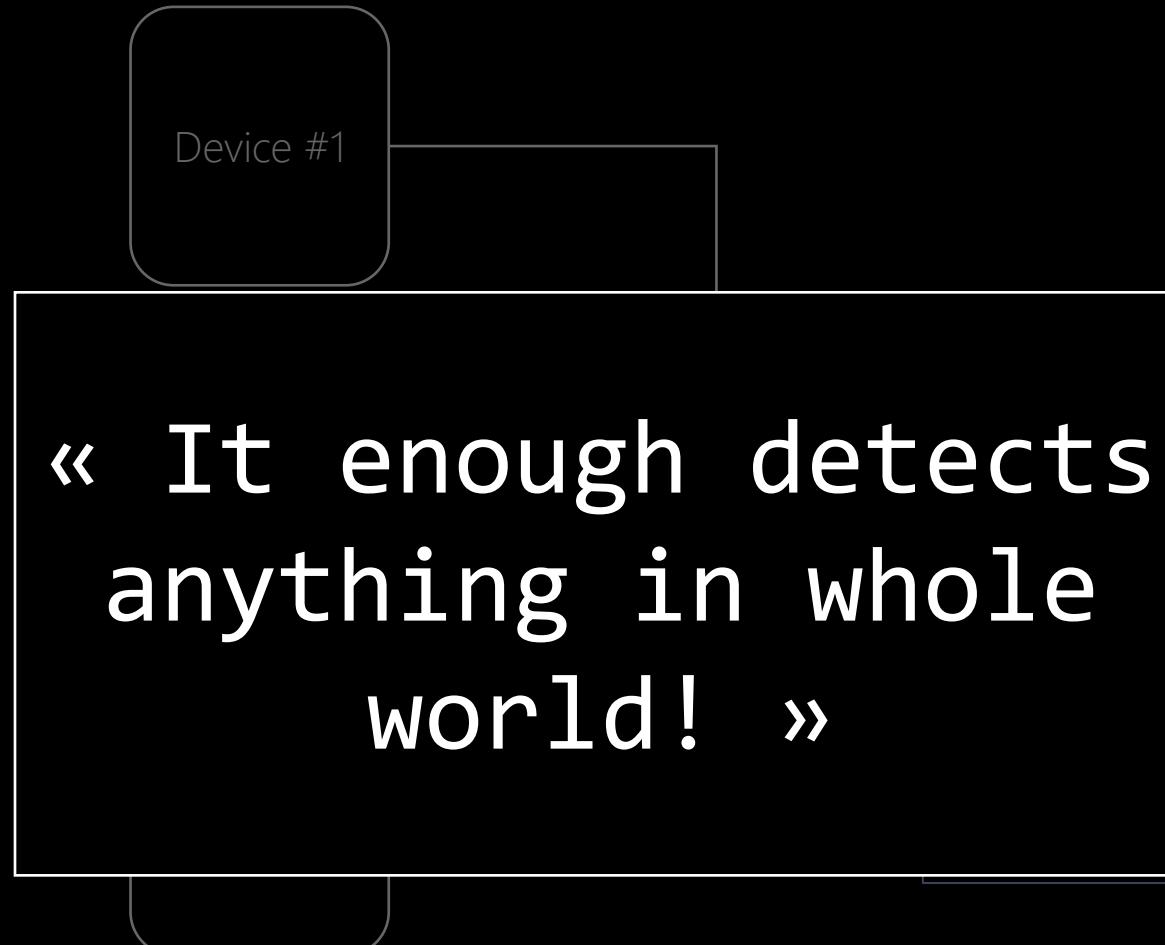
IIoT



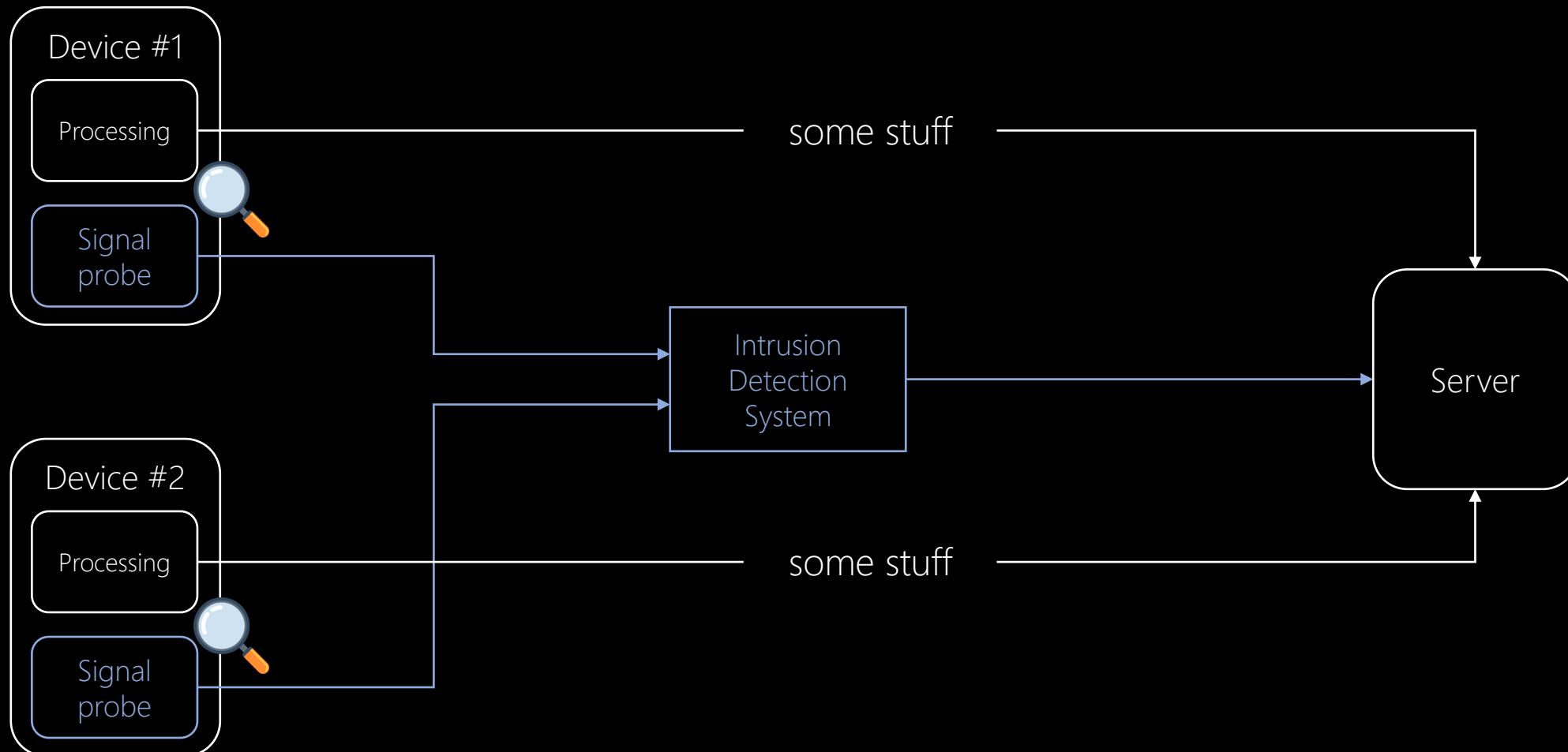
> What is an NIDS?



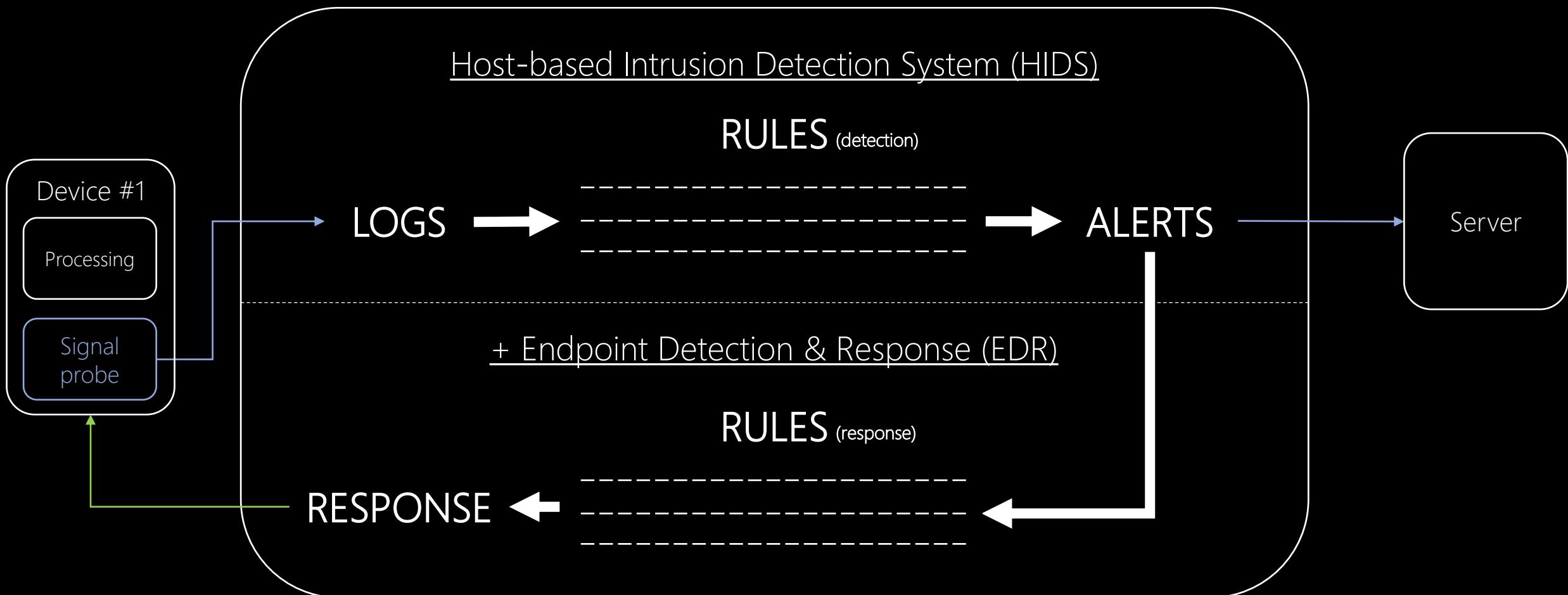
> What is an NIDS?



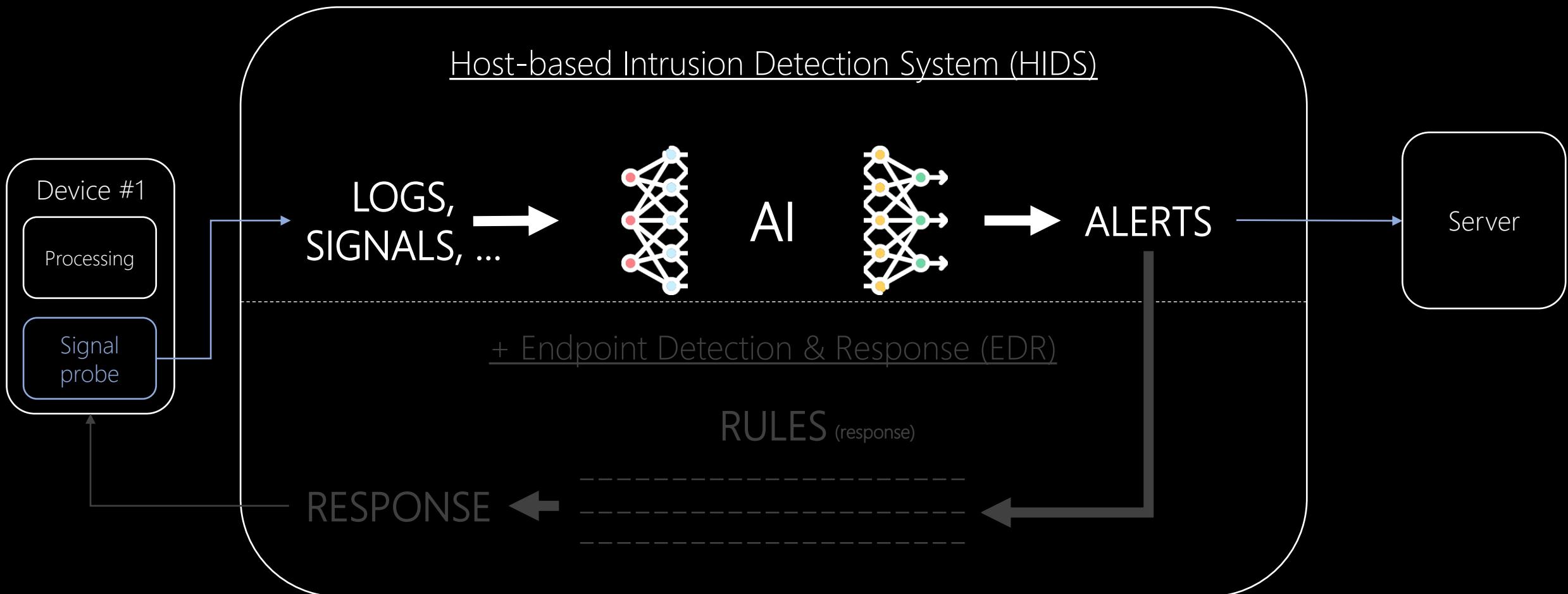
> What is an HIDS?



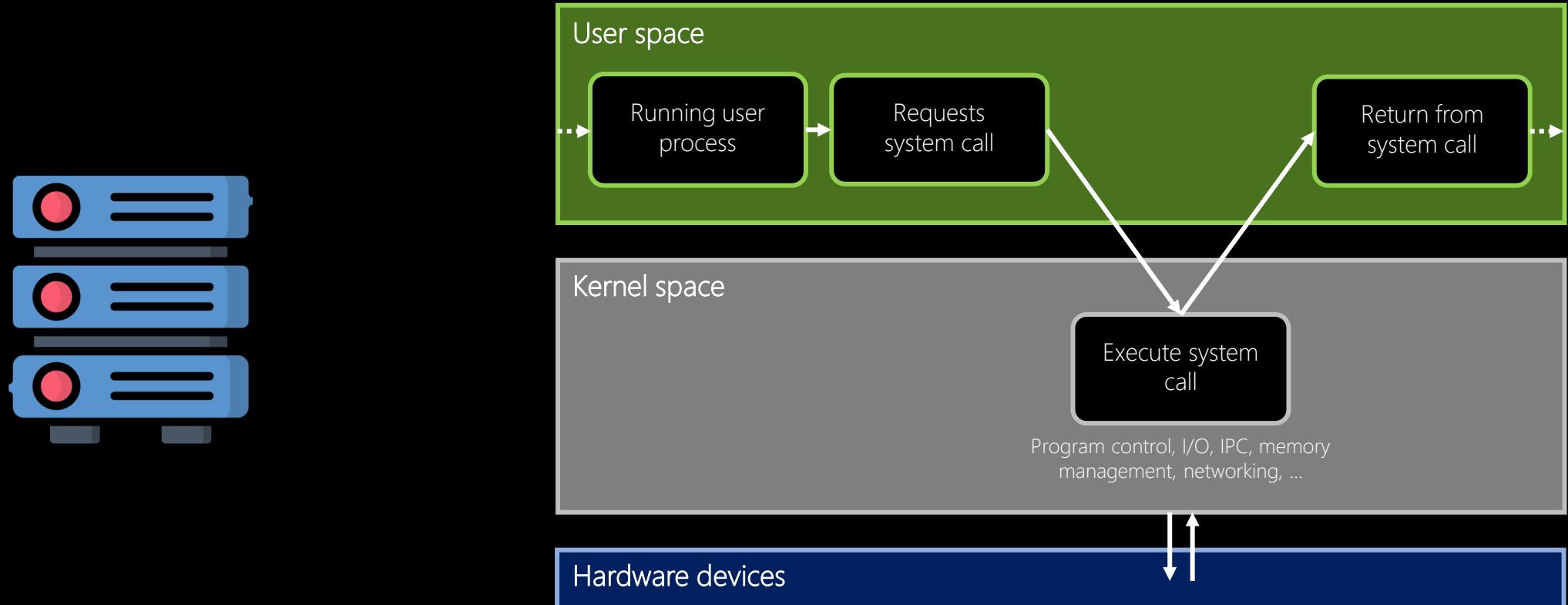
> What is an HIDS currently?



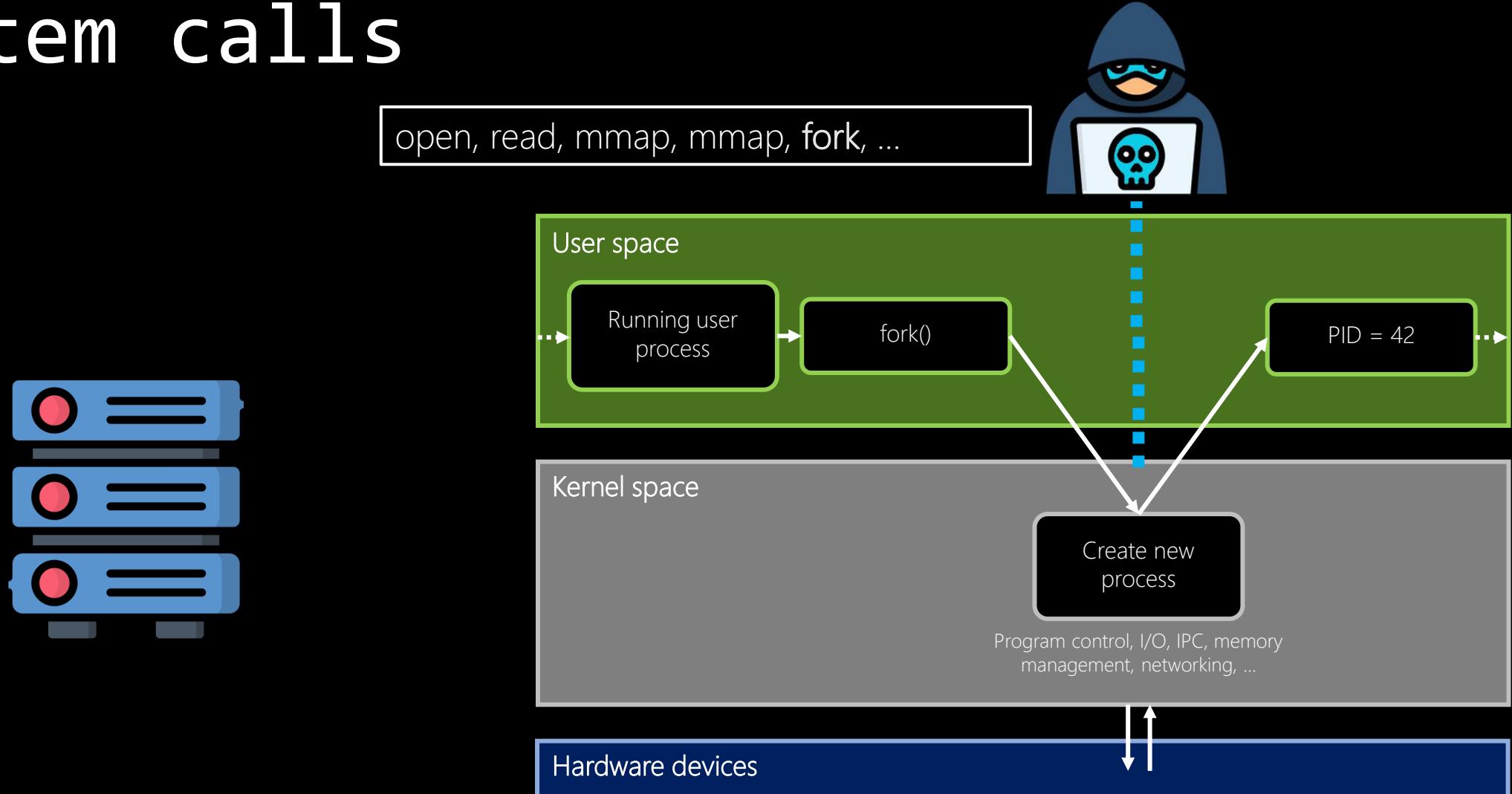
> What can be an HIDS soon?



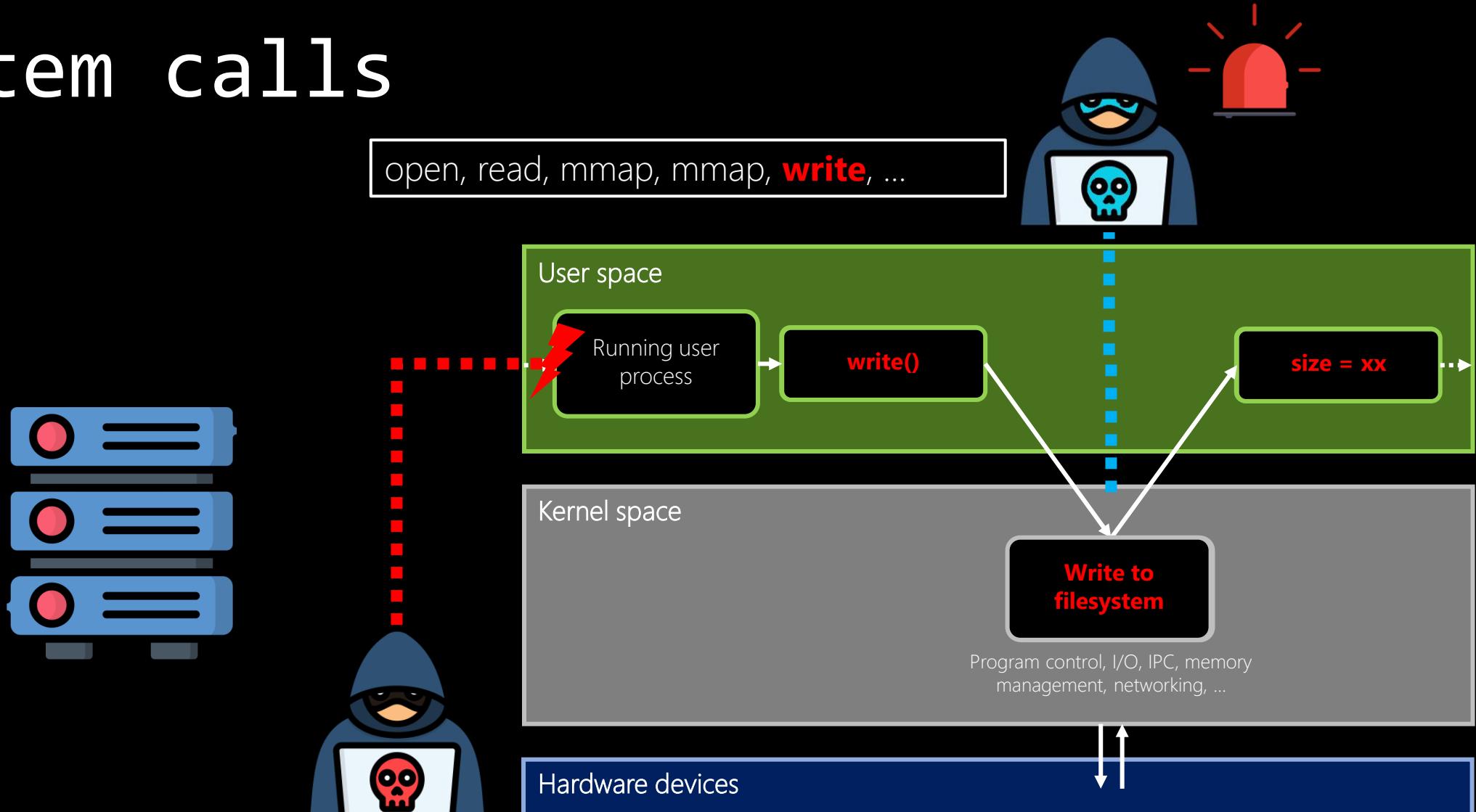
> System calls



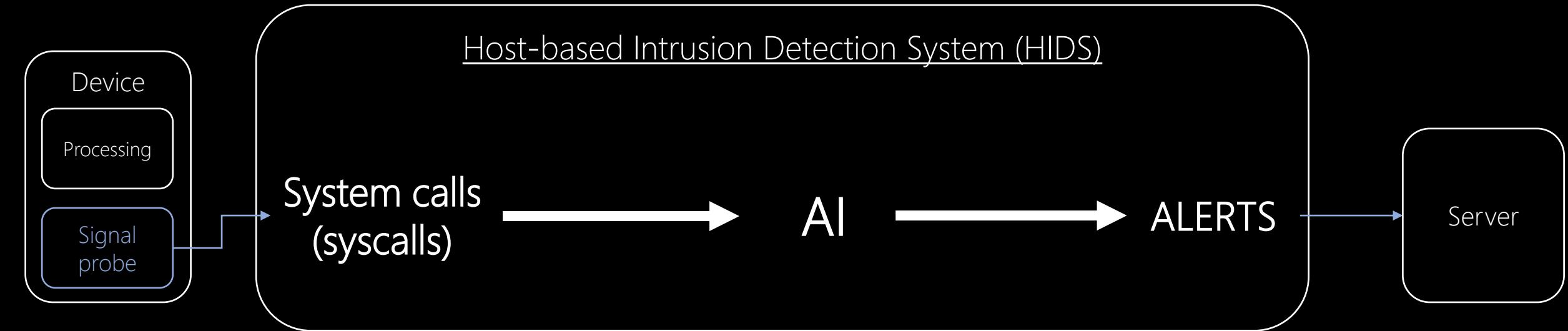
> System calls



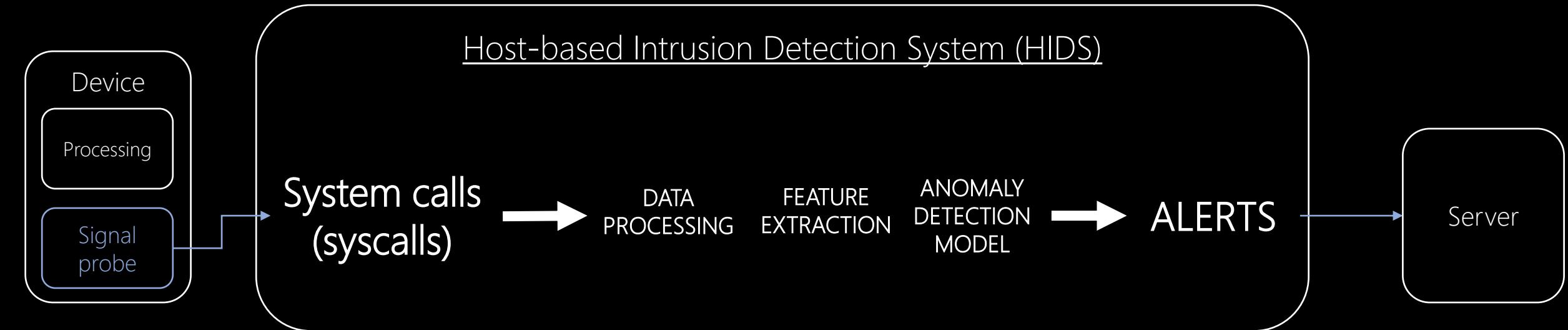
> System calls



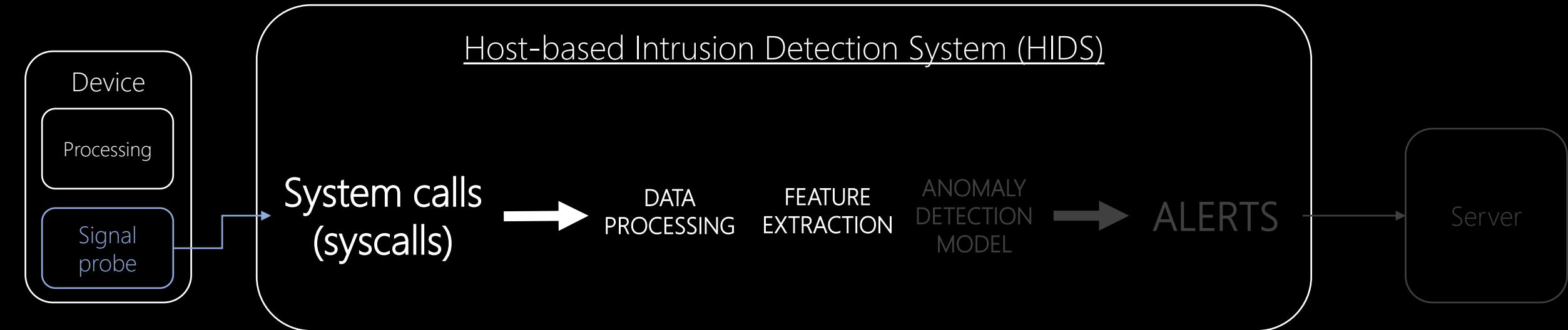
> What's the HIDS you'll deal with ?



> What's the HIDS you'll deal with ?



> What's the HIDS you'll deal with ?

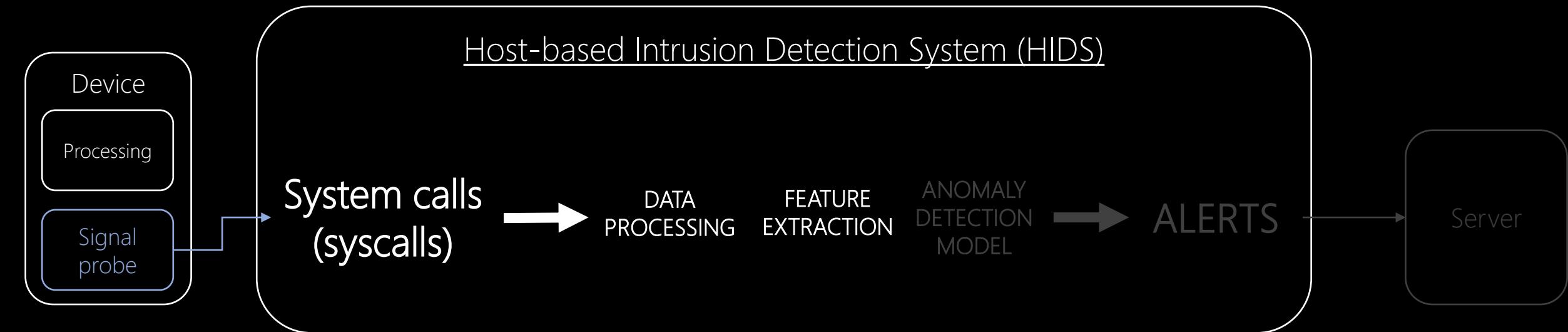


RAW SYSCALLS

Timestamp	ID	ARGS	RETVAL
00:01	#2	(2, 'hi')	2
00:02	#102	(3, {...})	3
00:04	#10	(39489)	&...
00:05	#102	{...}, {...}	{...}
00:07	#138	(409)	4
00:10

Trace #10

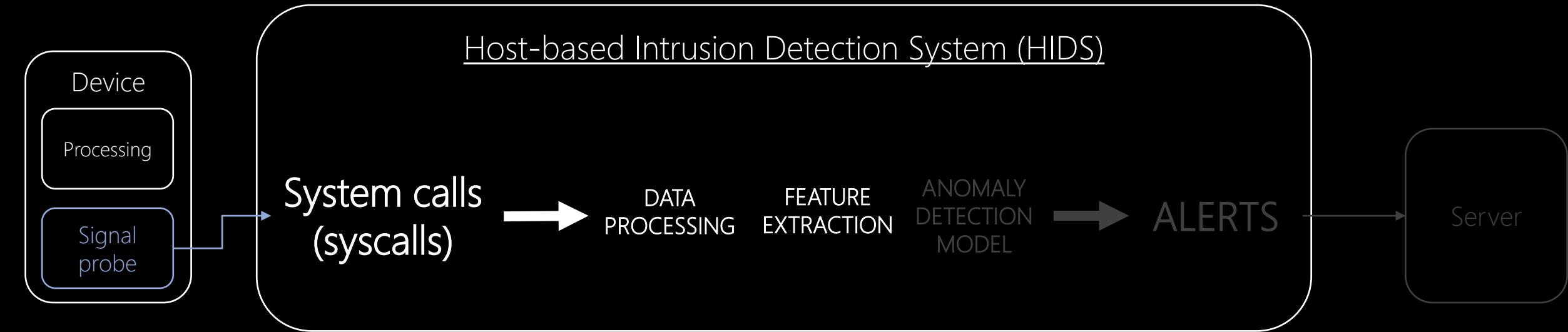
> What's the HIDS you'll deal with ?



Trace #10

RAW SYSCALLS				DATA PROCESSING	
Timestamp	ID	ARGS	RETVAL	ID	
00:01	#2	(2, 'hi')	2	#2	
00:02	#102	(3, {...})	3	#102	
00:04	#10	(39489)	&...	#10	
00:05	#102	(..., {...})	{...}	#102	
00:07	#138	(409)	4	#138	
00:10	

> What's the HIDS you'll deal with ?



Trace #10

RAW SYSCALLS			
Timestamp	ID	ARGS	RETVAL
00:01	#2	(2, 'hi')	2
00:02	#102	(3, {...})	3
00:04	#10	(39489)	&...
00:05	#102	(..., {...})	{...}
00:07	#138	(409)	4
00:10

DATA PROCESSING

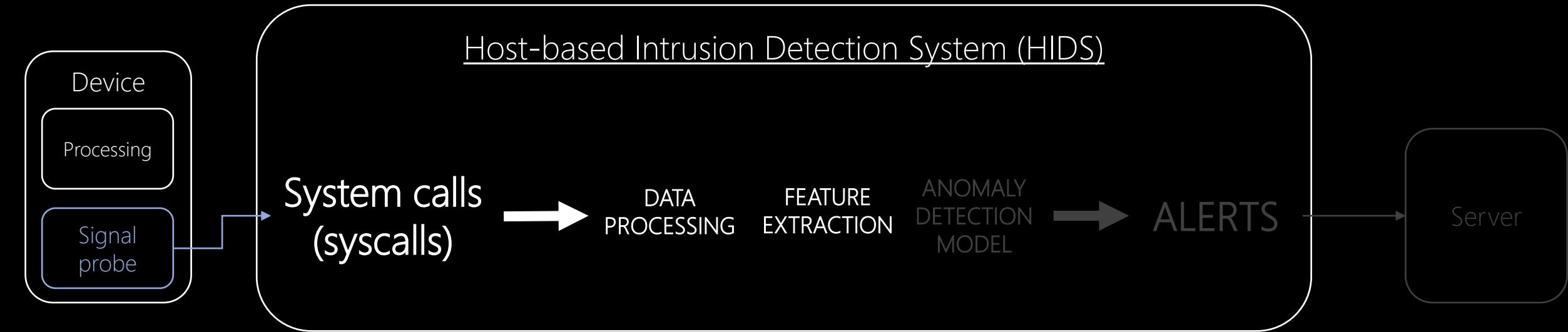
ID
#2
#102
#10
#102
#138
...

FEATURE EXTRACTION

3gram-ID
(#2, #102, #10)
(#102, #10, #102)
(#10, #102, #138)
(#102, #138, ...)
(#138, ..., ...)
...

A red box highlights the row for ID #102 in the DATA PROCESSING table, and a red arrow points from it to the corresponding 3gram-ID entry in the FEATURE EXTRACTION table.

> What's the HIDS you'll deal with ?



RAW SYSCALLS

Timestamp	ID	ARGS	RETVAL
00:01	#2	(2, 'hi')	2
00:02	#102	(3, {...})	3
00:04	#10	(39489)	&...
00:05	#102	{...}, {...}	{...}
00:07	#138	(409)	4
00:10

DATA PROCESSING

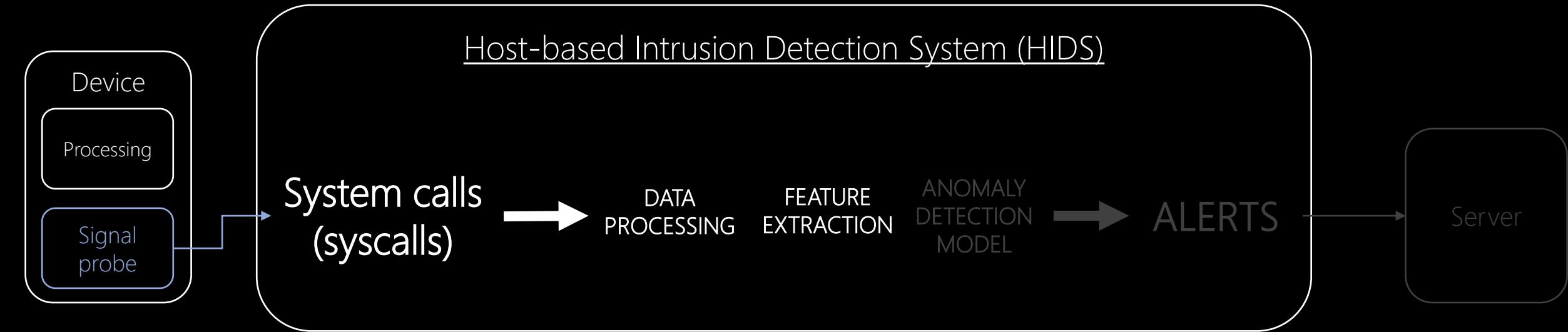
ID
#2
#102
#10
#102
#138
...

FEATURE EXTRACTION

3gram-ID
(#2, #102, #10)
(#102, #10, #102)
(#10, #102, #138)
(#102, #138, ...)
(#138,..., ...)
...

Trace #10

> What's the HIDS you'll deal with ?



RAW SYSCALLS

Timestamp	ID	ARGS	RETVAL
00:01	#2	(2, 'hi')	2
00:02	#102	(3, {...})	3
00:04	#10	(39489)	&...
00:05	#102	{...}, {...}	{...}
00:07	#138	(409)	4
00:10

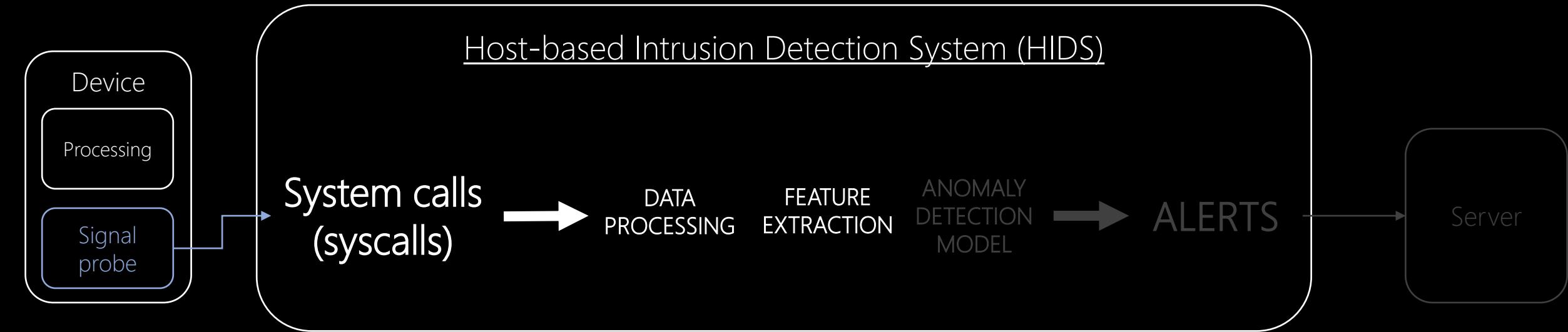
DATA PROCESSING

ID
#2
#102
#10
#102
#138
...

FEATURE EXTRACTION

3gram-ID
(#2, #102, #10)
(#102, #10, #102)
(#10, #102, #138)
(#102, #138, ...)
(#138,..., ...)
...

> What's the HIDS you'll deal with ?



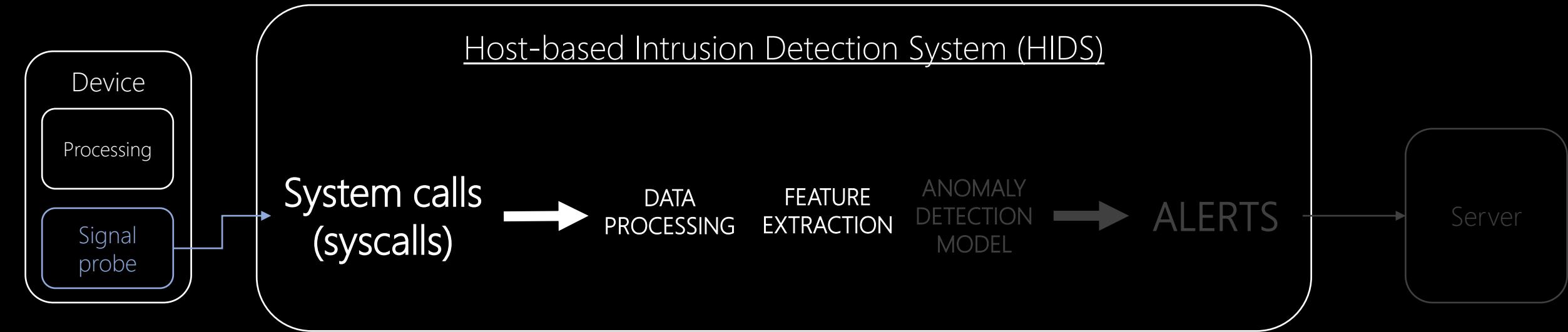
Trace #10

RAW SYSCALLS			
Timestamp	ID	ARGS	RETVAL
00:01	#2	(2, 'hi')	2
00:02	#102	(3, {...})	3
00:04	#10	(39489)	&...
00:05	#102	{...}, {...}	{...}
00:07	#138	(409)	4
00:10

DATA PROCESSING	
ID	3gram-ID
#2	(#2, #102, #10)
#102	(#102, #10, #102)
#10	(#10, #102, #138)
#102	(#102, #138, ...)
#138	(#138, ..., ...)
...	...

FEATURE EXTRACTION	
--------------------	--

> What's the HIDS you'll deal with ?



Trace #10

RAW SYSCALLS			
Timestamp	ID	ARGS	RETVAL
00:01	#2	(2, 'hi')	2
00:02	#102	(3, {...})	3
00:04	#10	(39489)	&...
00:05	#102	{...}, {...}	{...}
00:07	#138	(409)	4
00:10

DATA PROCESSING

ID
#2
#102
#10
#102
#138
...

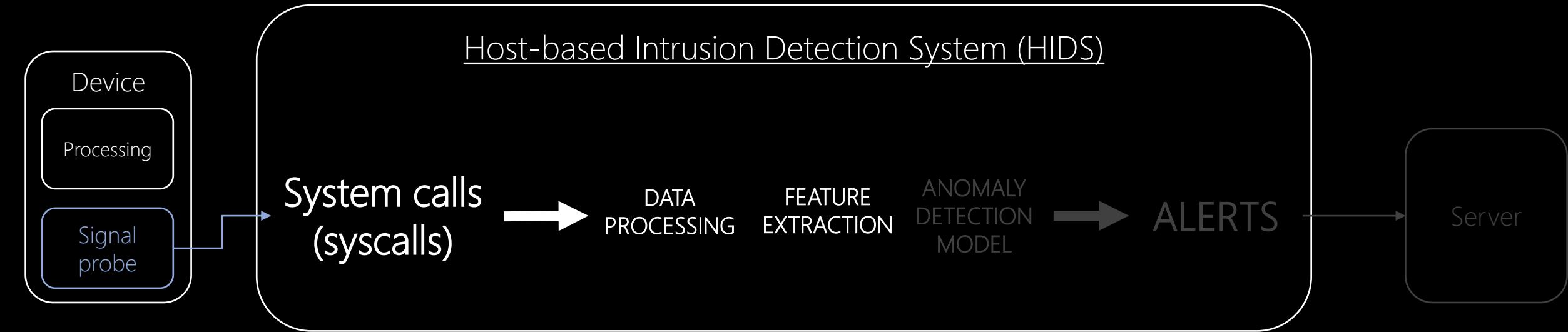
FEATURE EXTRACTION

$TF-IDF = N\text{-gram frequency in trace} * \log(\frac{\text{number of past traces}}{\text{number of past traces with } N\text{-gram}})$

3gram-ID
(#2, #102, #10)
(#102, #10, #102)
(#10, #102, #138)
(#102, #138, ...)
(#138, ..., ...)
...

TF-IDF
0.3
0.5
0.02
0.9
0.4
...

> What's the HIDS you'll deal with ?



Trace #10

RAW SYSCALLS			
Timestamp	ID	ARGS	RETVAL
00:01	#2	(2, 'hi')	2
00:02	#102	(3, {...})	3
00:04	#10	(39489)	&...
00:05	#102	{...}, {...}	{...}
00:07	#138	(409)	4
00:10

DATA PROCESSING

ID
#2
#102
#10
#102
#138
...

FEATURE EXTRACTION

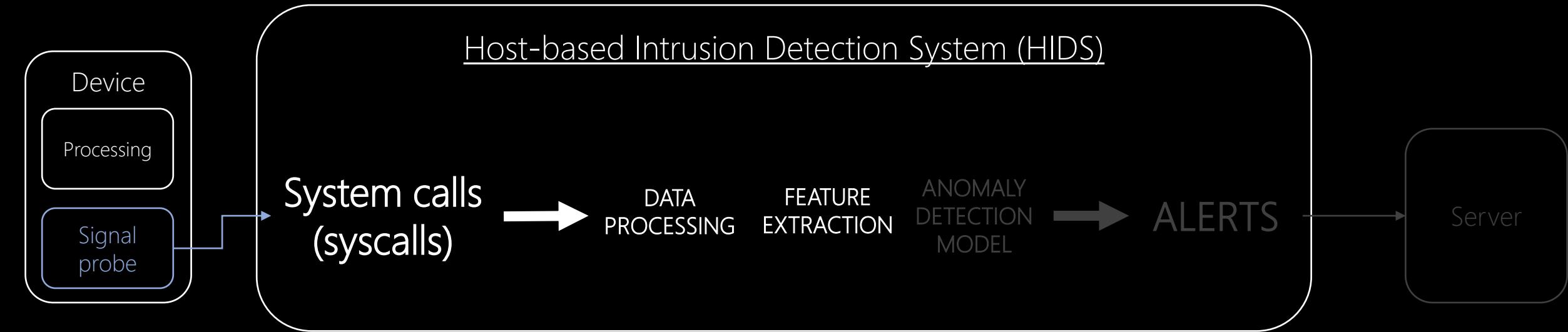
$TF-IDF = N\text{-gram frequency in trace} * \log(\frac{\text{number of past traces}}{\text{number of past traces with } N\text{-gram}})$

3gram-ID
(#2, #102, #10)
(#102, #10, #102)
(#10, #102, #138)
(#102, #138, ...)
(#138, ..., ...)
...

TF-IDF
0.3
0.5
0.02
0.9
0.4
...

A red box highlights the entry '(#102, #10, #102)' in the 3gram-ID table, and a red arrow points from it to the corresponding TF-IDF value of 0.5.

> What's the HIDS you'll deal with ?



RAW SYSCALLS

Timestamp	ID	ARGS	RETVAL
00:01	#2	(2, 'hi')	2
00:02	#102	(3, {...})	3
00:04	#10	(39489)	&...
00:05	#102	{...}, {...}	{...}
00:07	#138	(409)	4
00:10

DATA PROCESSING

ID
#2
#102
#10
#102
#138
...

FEATURE EXTRACTION

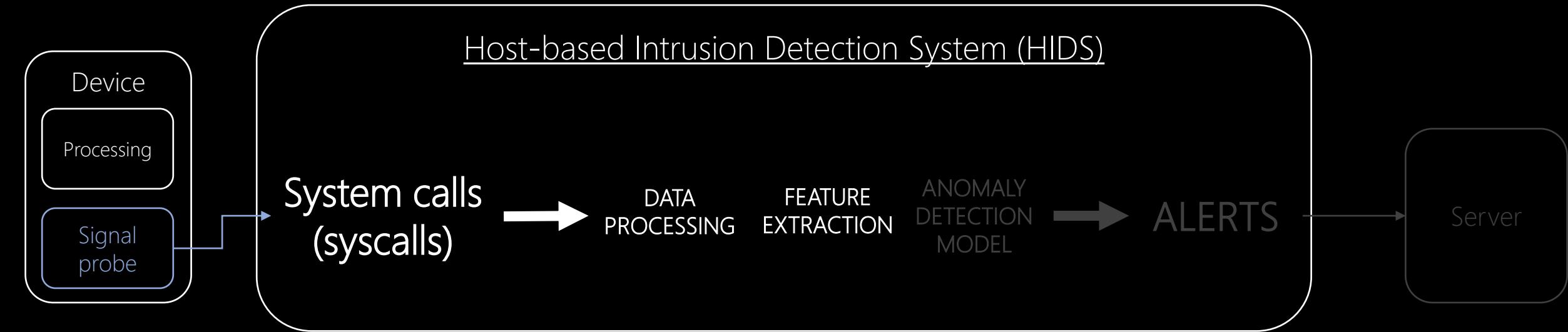
$TF-IDF = N\text{-gram frequency in trace} * \log(\frac{\text{number of past traces}}{\text{number of past traces with } N\text{-gram}})$

3gram-ID
(#2, #102, #10)
(#102, #10, #102)
(#10, #102, #138)
(#102, #138, ...)
(#138, ..., ...)
...

TF-IDF
0.3
0.5
0.02
0.9
0.4
...

Trace #10

> What's the HIDS you'll deal with ?



Trace #10

RAW SYSCALLS			
Timestamp	ID	ARGS	RETVAL
00:01	#2	(2, 'hi')	2
00:02	#102	(3, {...})	3
00:04	#10	(39489)	&...
00:05	#102	{...}, {...}	{...}
00:07	#138	(409)	4
00:10

DATA PROCESSING

ID
#2
#102
#10
#102
#138
...

FEATURE EXTRACTION

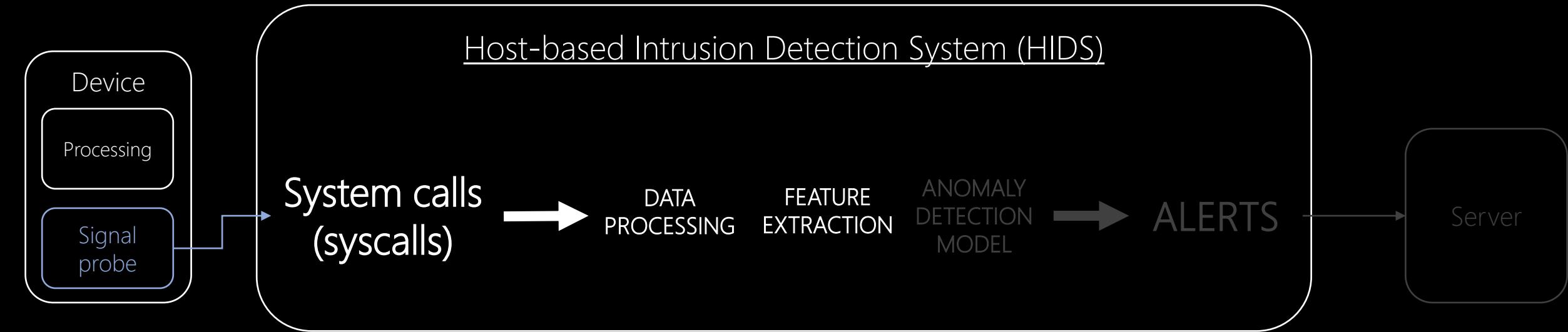
$TF-IDF = N\text{-gram frequency in trace} * \log(\frac{\text{number of past traces}}{\text{number of past traces with } N\text{-gram}})$

3gram-ID
(#2, #102, #10)
(#102, #10, #102)
(#10, #102, #138)
(#102, #138, ...)
(#138, ...)
...

TF-IDF
0.3
0.5
0.02
0.9
0.4
...

A red box highlights the row '(#102, #138, ...)'. A red arrow points from this row to the corresponding TF-IDF value of 0.9.

> What's the HIDS you'll deal with ?



Trace #10

Timestamp	ID	ARGS	RETVAL
00:01	#2	(2, 'hi')	2
00:02	#102	(3, {...})	3
00:04	#10	(39489)	&...
00:05	#102	{...}, {...}	{...}
00:07	#138	(409)	4
00:10

DATA PROCESSING

ID
#2
#102
#10
#102
#138
...

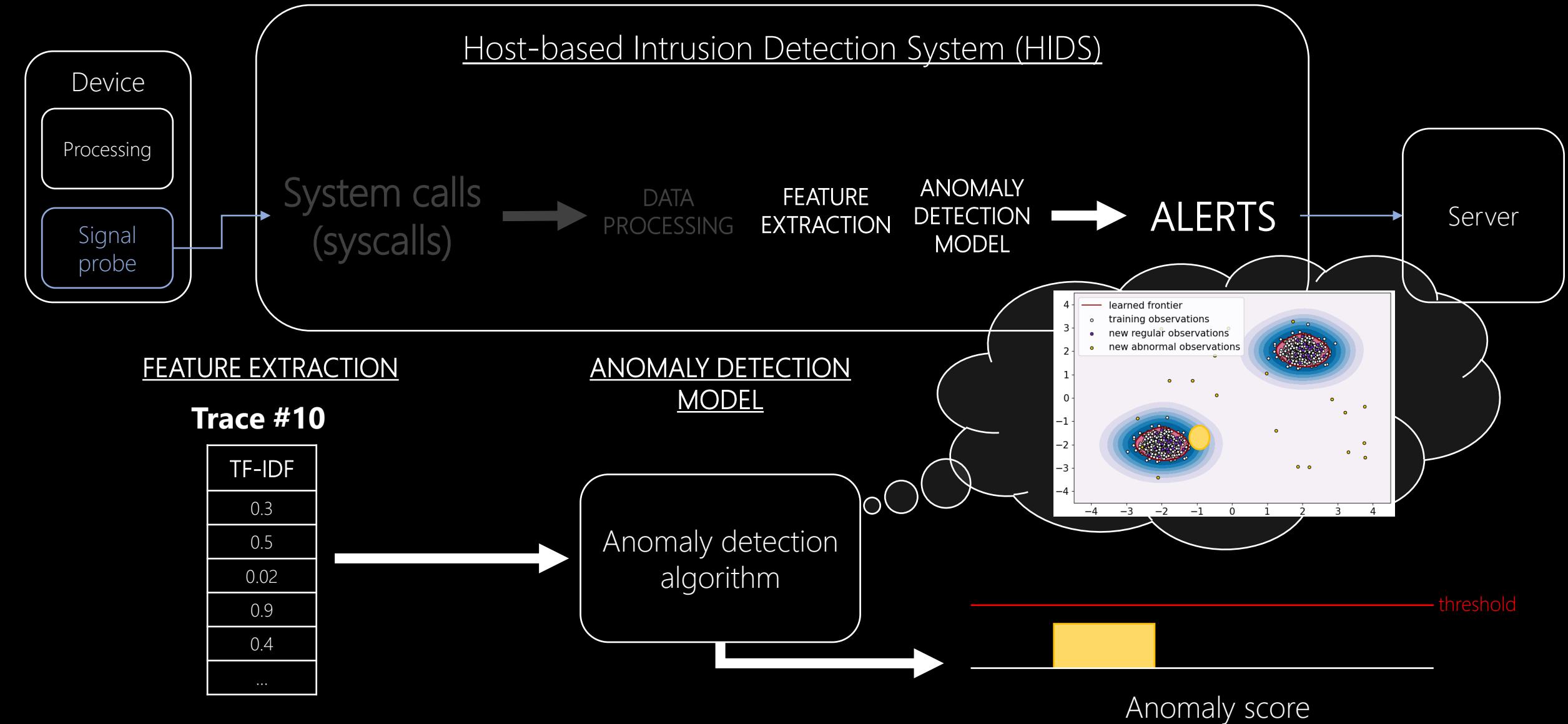
FEATURE EXTRACTION

$TF-IDF = N\text{-gram frequency in trace} * \log(\frac{\text{number of past traces}}{\text{number of past traces with } N\text{-gram}})$

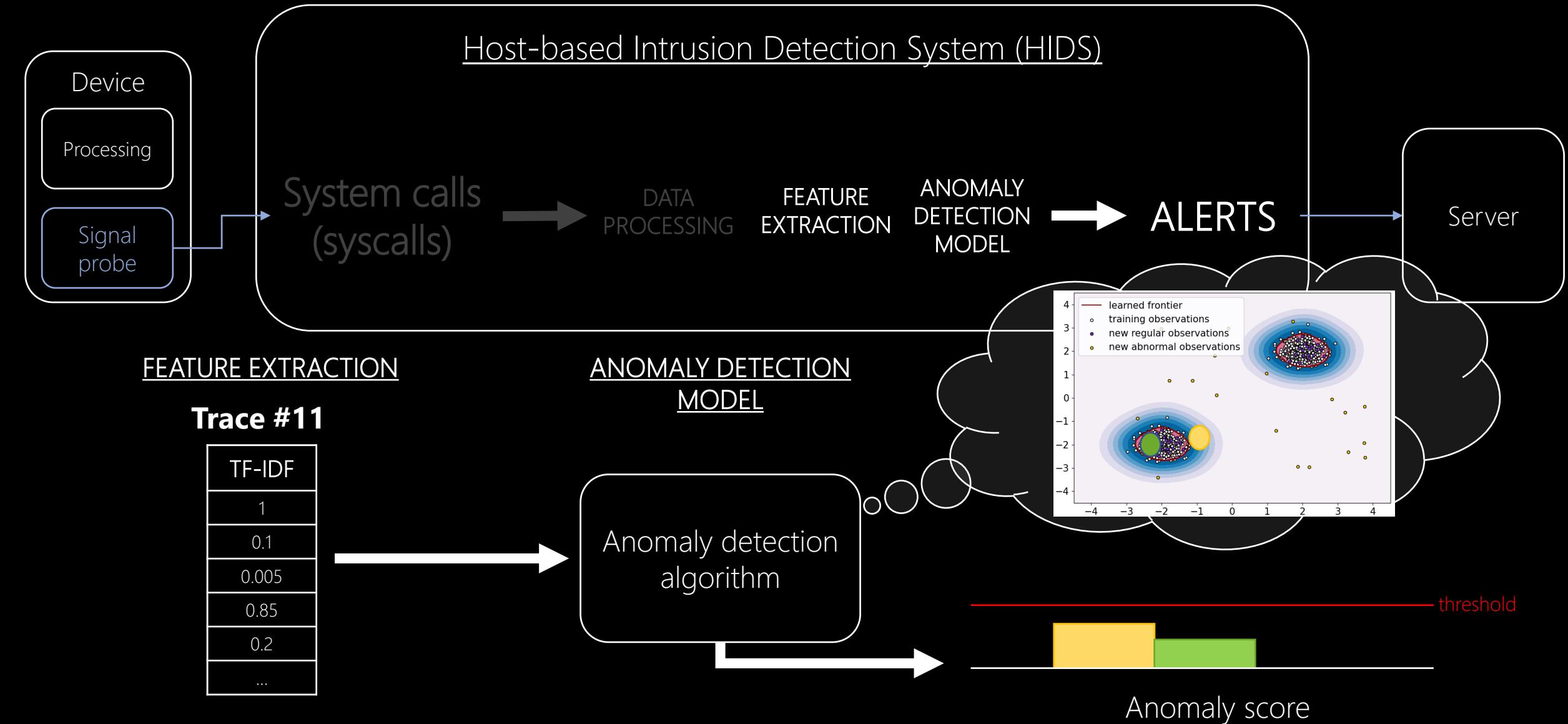
3gram-ID
(#2, #102, #10)
(#102, #10, #102)
(#10, #102, #138)
(#102, #138, ...)
(#138, ..., ...)
...

TF-IDF
0.3
0.5
0.02
0.9
0.4
...

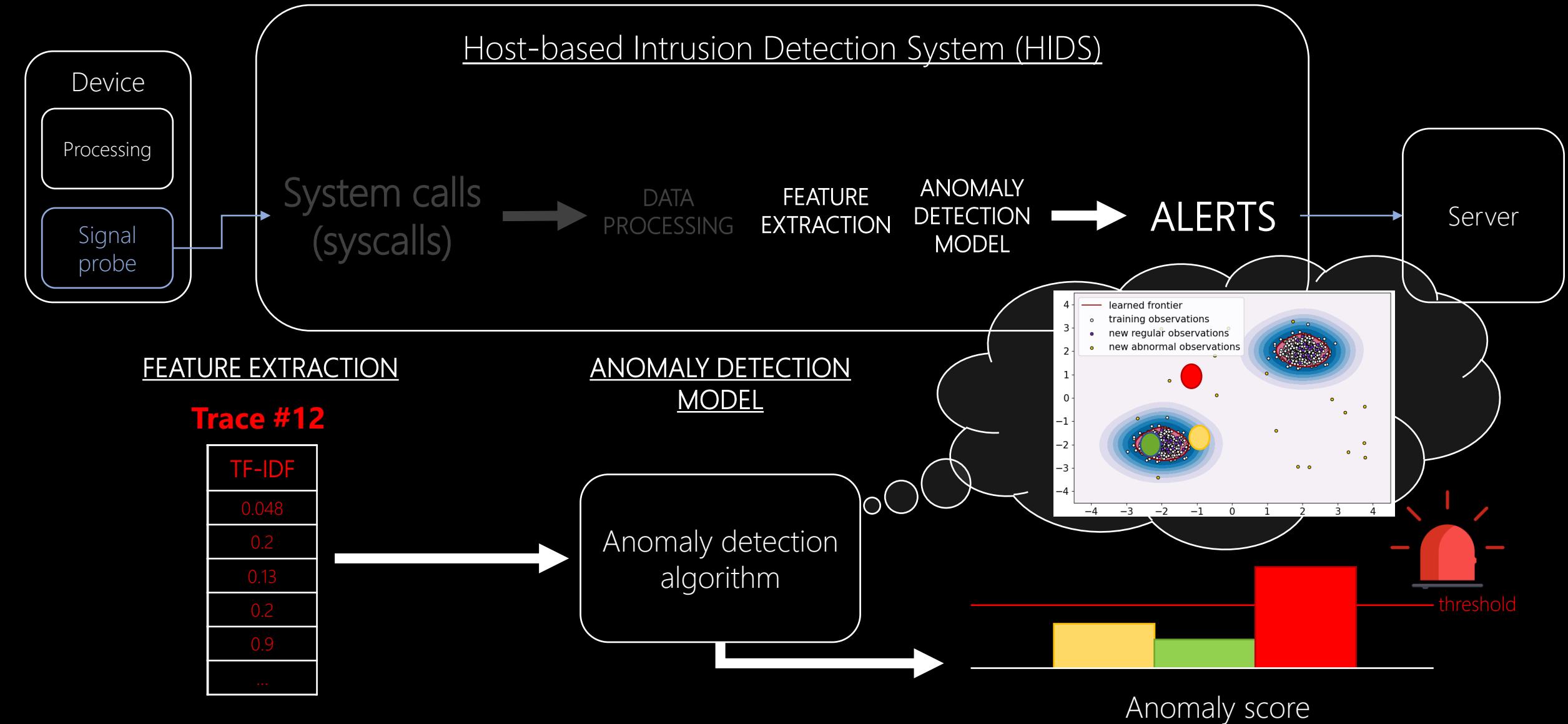
> What's the HIDS you'll deal with ?



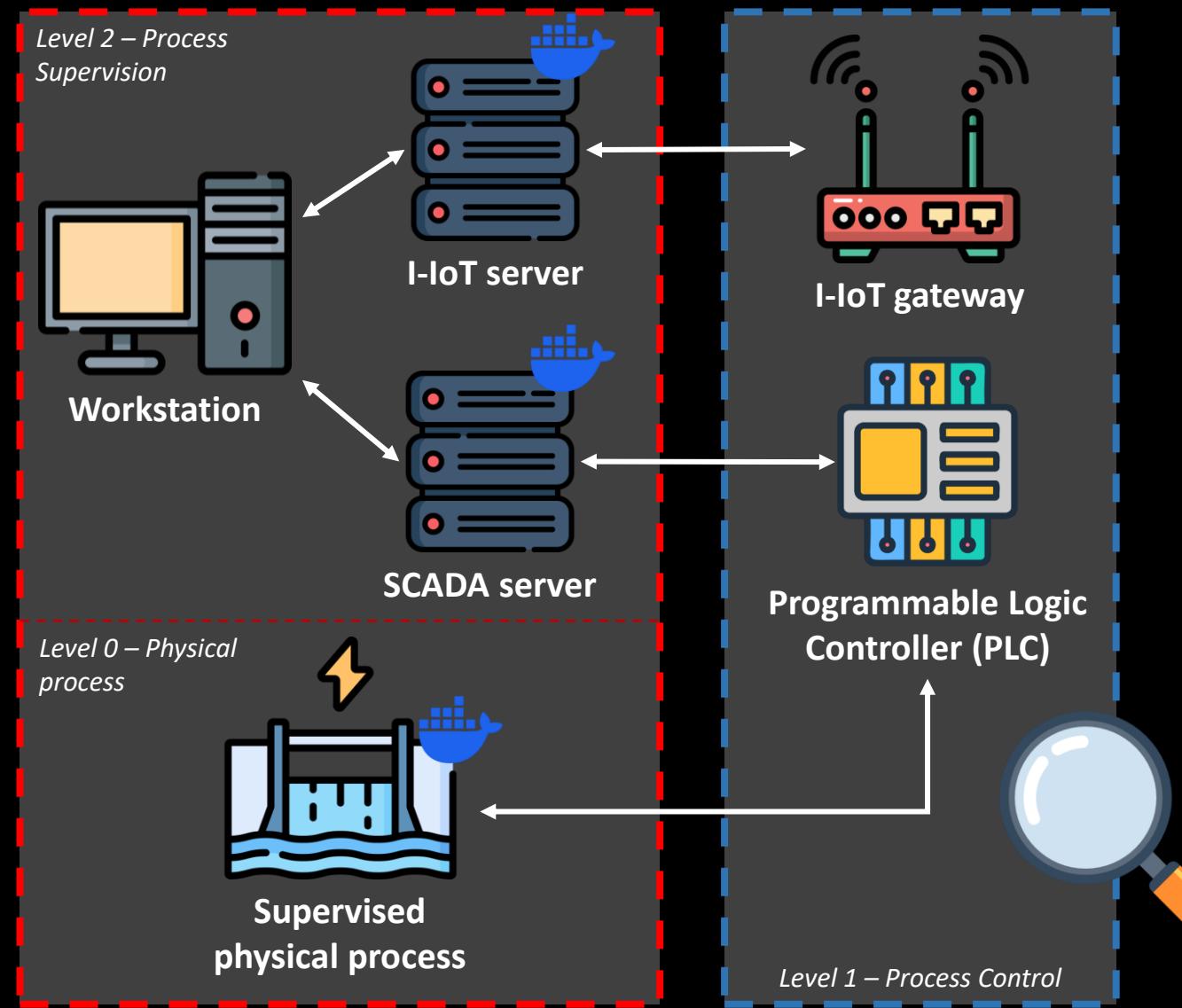
> What's the HIDS you'll deal with ?



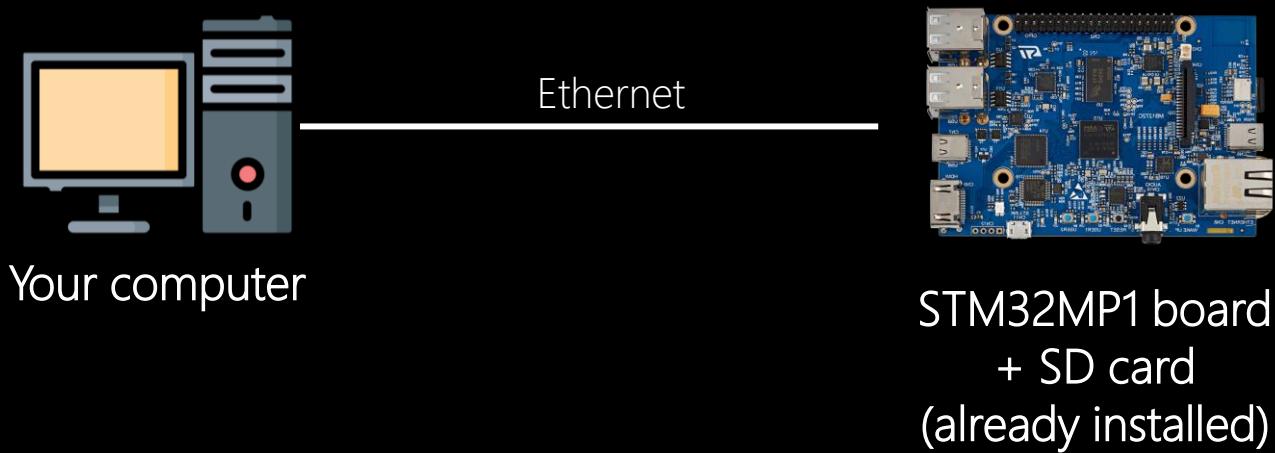
> What's the HIDS you'll deal with ?



> Hands-on: use case

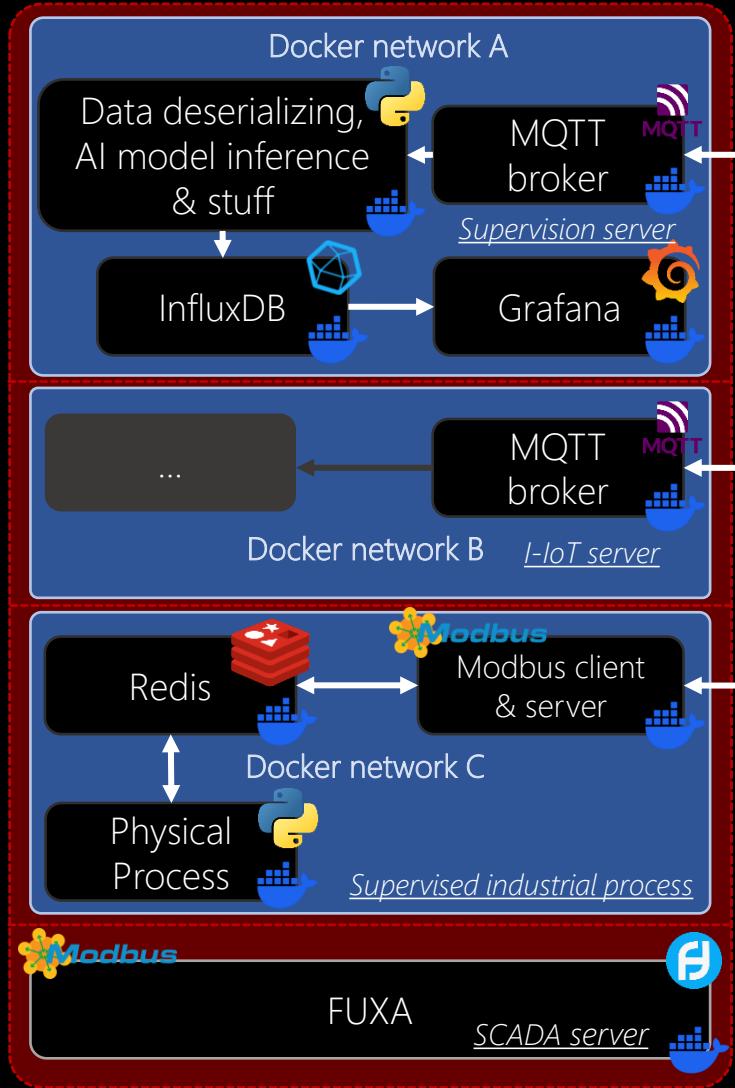


> Hands-on: platform (hardware)

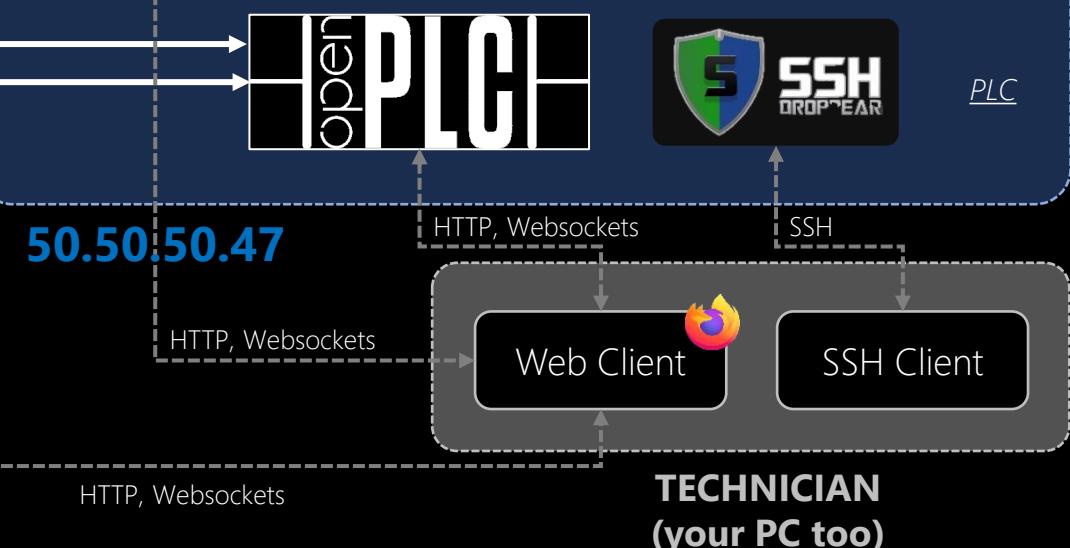
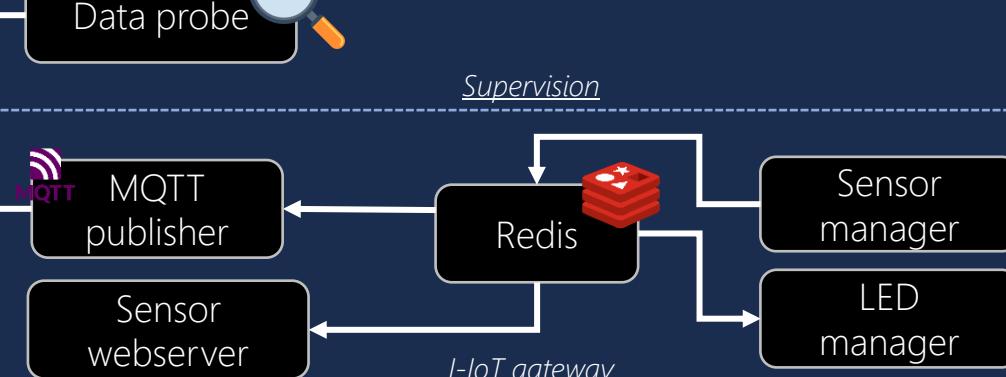


> Hands-on: platform (running services)

SERVER (your PC)



INDUSTRIAL DEVICE (STM32MP1)



> Hands-on: threat model

An attacker (you) is in the same network as the PLC

(this is a critical scenario for a real-life facility, but it has already been seen multiple times...)

WHAT ARE THE RULES?

1. Suppose you do not have physical access to the device
2. Do not try to attack things on your computer (industrial process simulator, AI & supervision stuff); that's not the goal of the workshop.
- X. There are hidden flags on the target... Try to find 'em, without being caught by the IDS ☺

> Installation + Demo

30min

Everything that needs to run on your computer: **github.com/LalieA/Workshop_AI_Protected_Realm**



For further information after GreHack, don't hesitate to reach us:
lalie.arnoud@cea.fr
ulysse.vincenti@cea.fr

> Now, it's up to you!

45min



> Closing remarks

This intrusion detection system is flawed (mimicry attacks)

```
open('file.txt', O_RDWR) = 7  
mmap(0xae7fd40b, 256, ...)  
mmap(0xfe5a3b16, 512, ...)  
read(7, &string_var, 20)  
    fork()  
write(7, &string_var_2, 20)
```

Recorded good behavior,
learned by anomaly
detection model

```
open('/etc/passwd', O_RDWR) = 10  
Fails → mmap(0xffffffff, 0, ...)  
Fails → mmap(0xffffffff, 0, ...)  
Useless → read(10, &string_var, 1024)  
        fork()  
        write(10, &new_line_for_user_var, 512)  
        ...
```

Sequence that is learned
as valid, but hides the
attacker's behavior

> Closing remarks

AI helps solving a problem, but brings others

- It is difficult to explain anomaly scores (see **Explainable AI**)
- Need other mechanisms for forensics after an alert is raised
- AI models can be resource consuming and hardly integratable as is, especially in embedded systems (see **Pruning, Quantization, Distillation**)
- AI models are vulnerable to **data poisoning** and **adversarial attacks**
- **LLMs are not magic** 🤖 🧙

> Closing remarks

This EDR, as others, is still vulnerable to EDR evasion techniques

- Payloads can **behave too similarly** to known benign processes
- System calls' **hooking mechanism** can be altered
- **Indirect system calls** can be made

A 6-months internship is available in
CEA Grenoble on bluetooth fuzzing
→ ulysse.vincenti@cea.fr

A 6-months internship is available in
CEA Grenoble on the hardening of an
EDR for embedded systems, similar to
what has been seen in this workshop.
→ lalie.arnoud@cea.fr

> Thank you for participating
in this workshop!

If you have feedbacks, questions, or just want to drink a beer, don't hesitate to reach us:

lalie.arnoud@cea.fr

ulysses.vincenti@cea.fr