

CYBERSECURITY AIs

From PentestGPT to Building an
AI-Powered Robot Immune System

VÍCTOR MAYORAL VILCHES

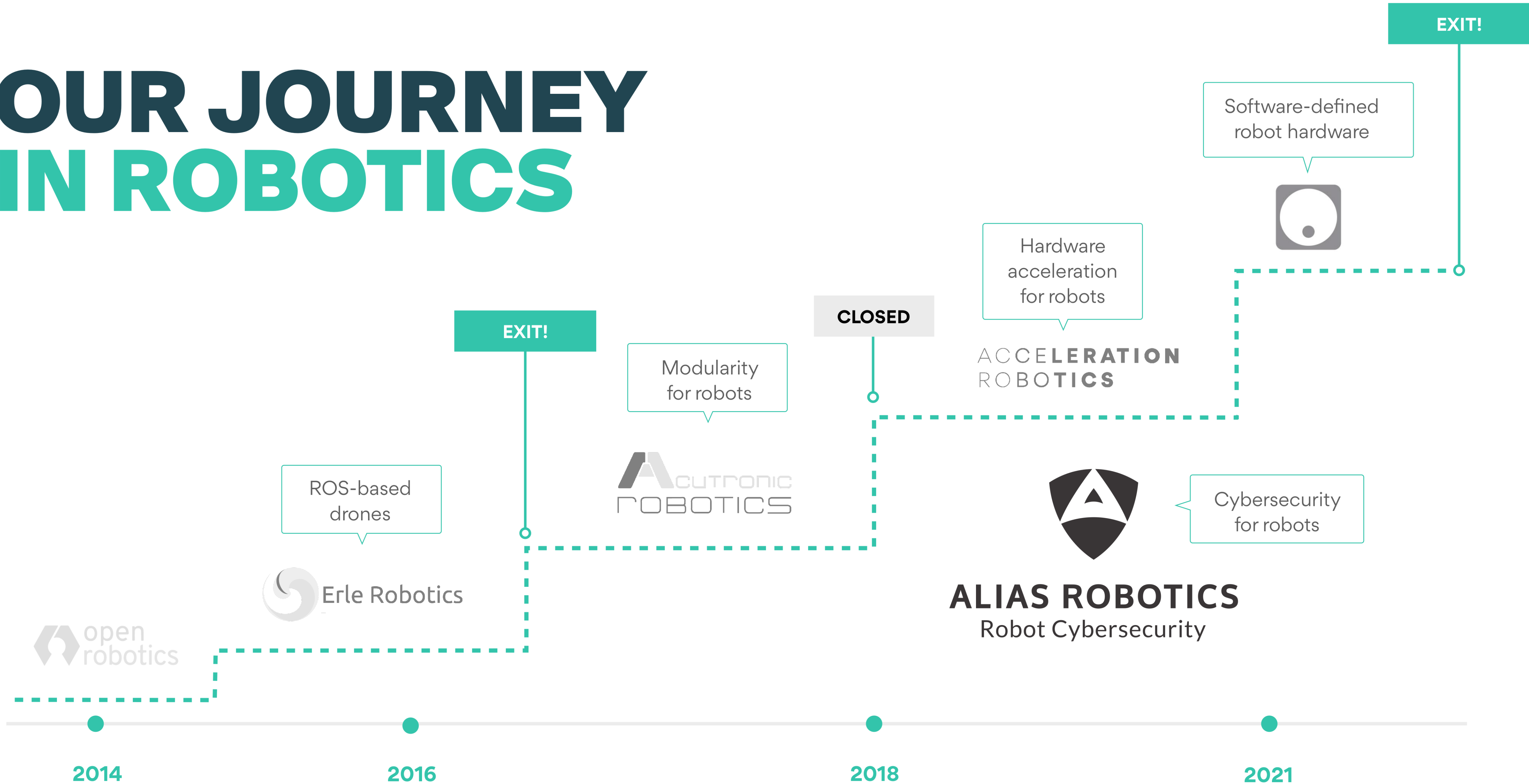
Founder

victor@aliasrobotics.com



ALIAS ROBOTICS
Robot Cybersecurity

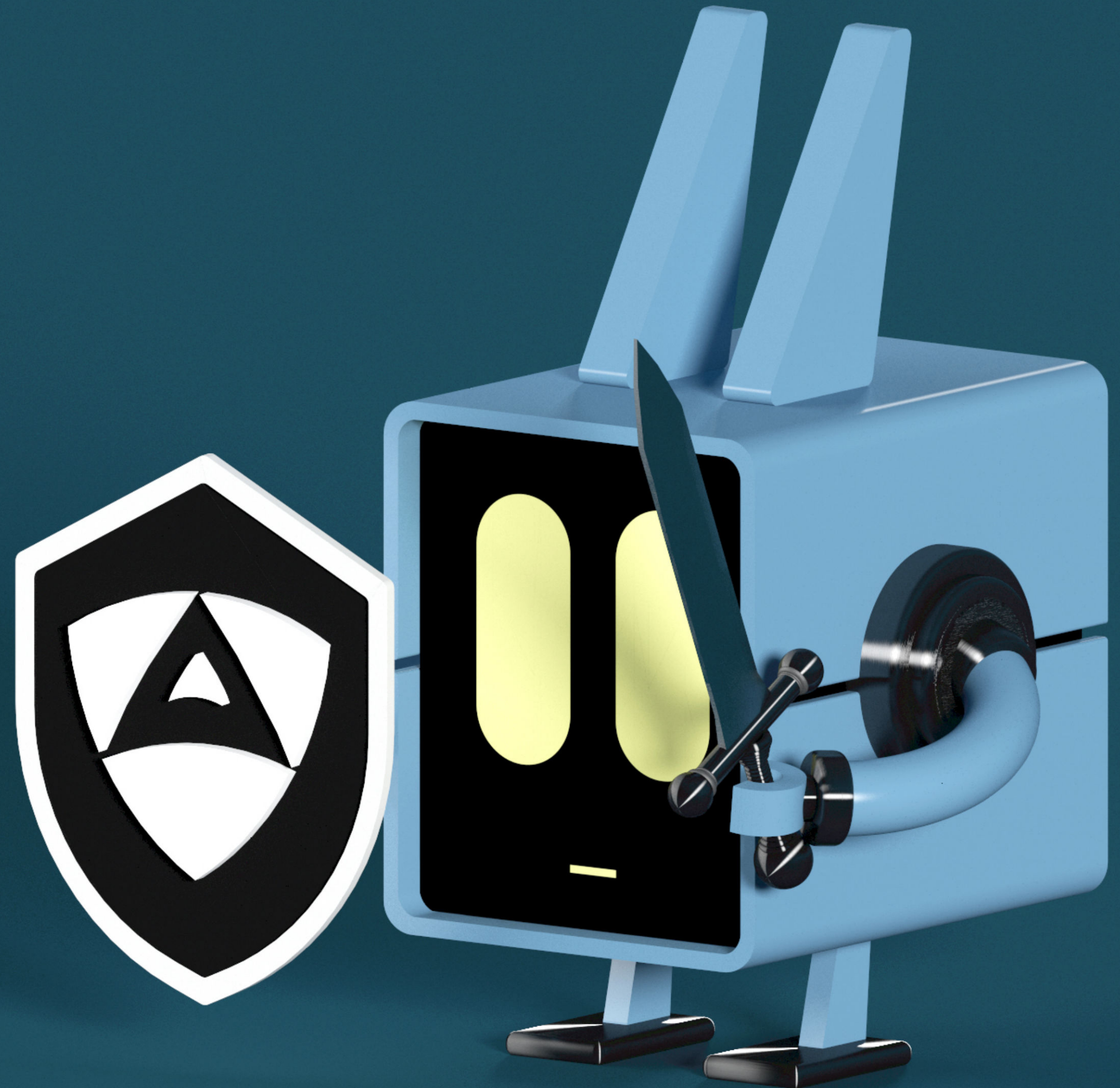
OUR JOURNEY IN ROBOTICS



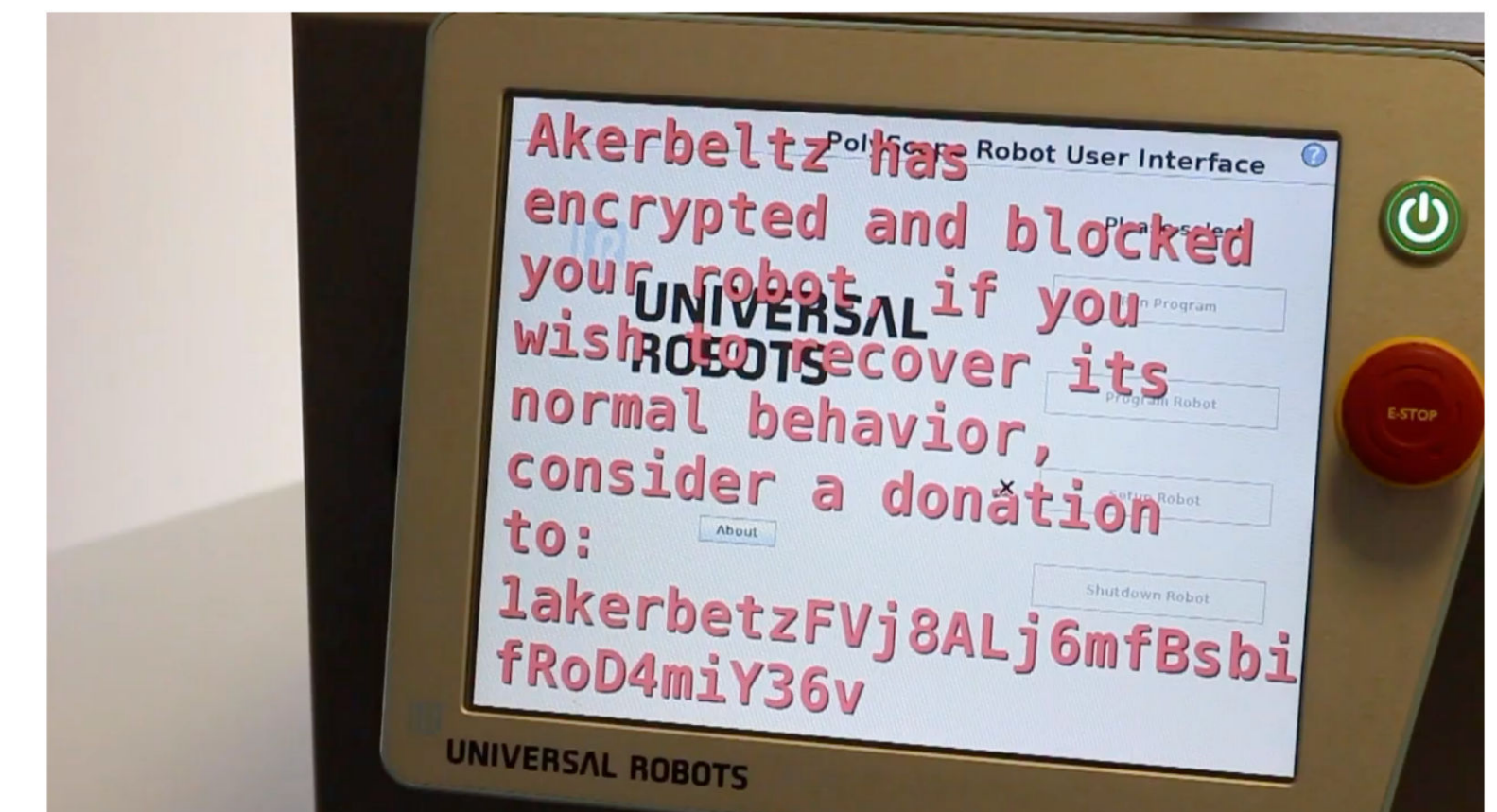
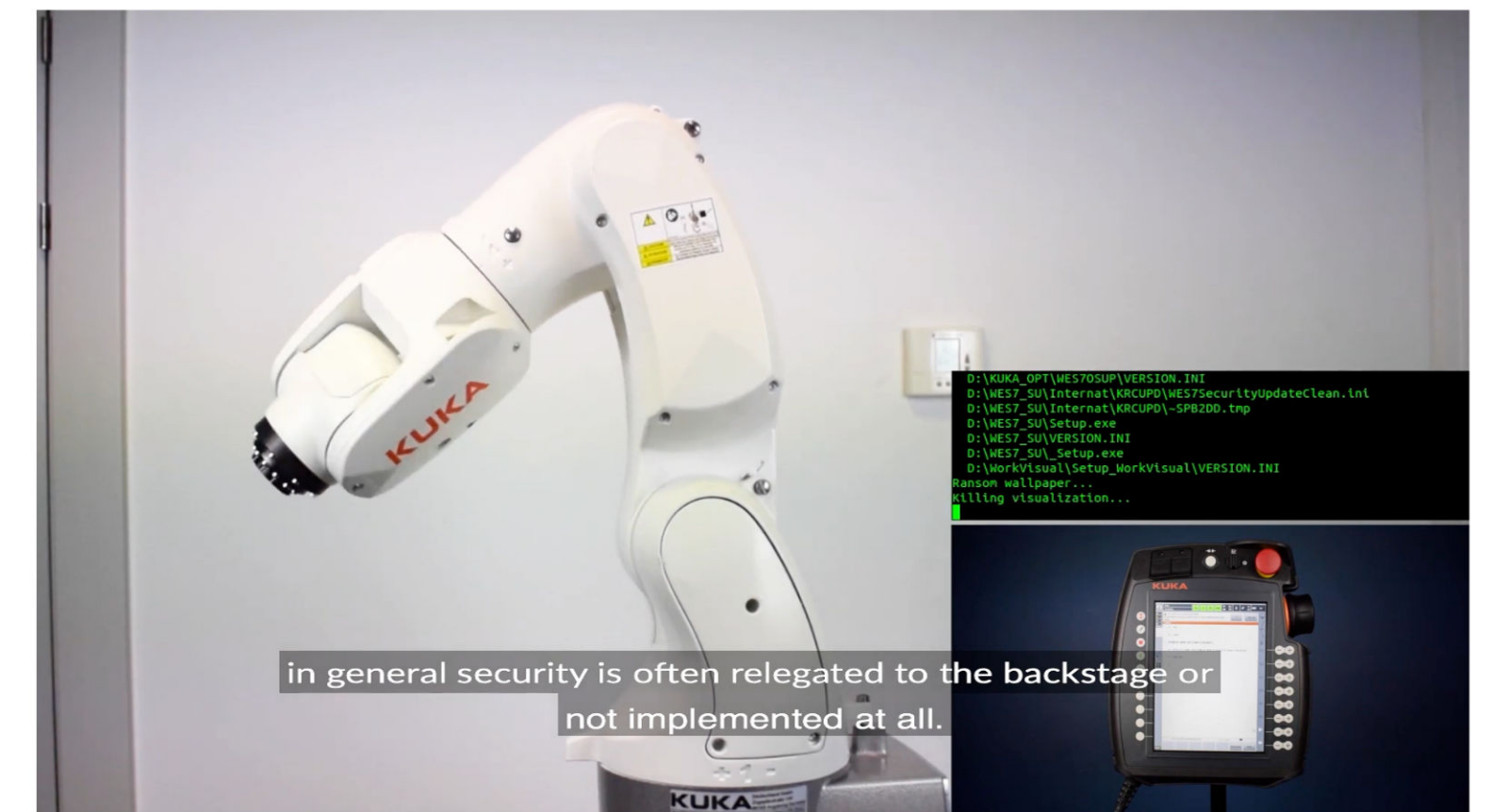
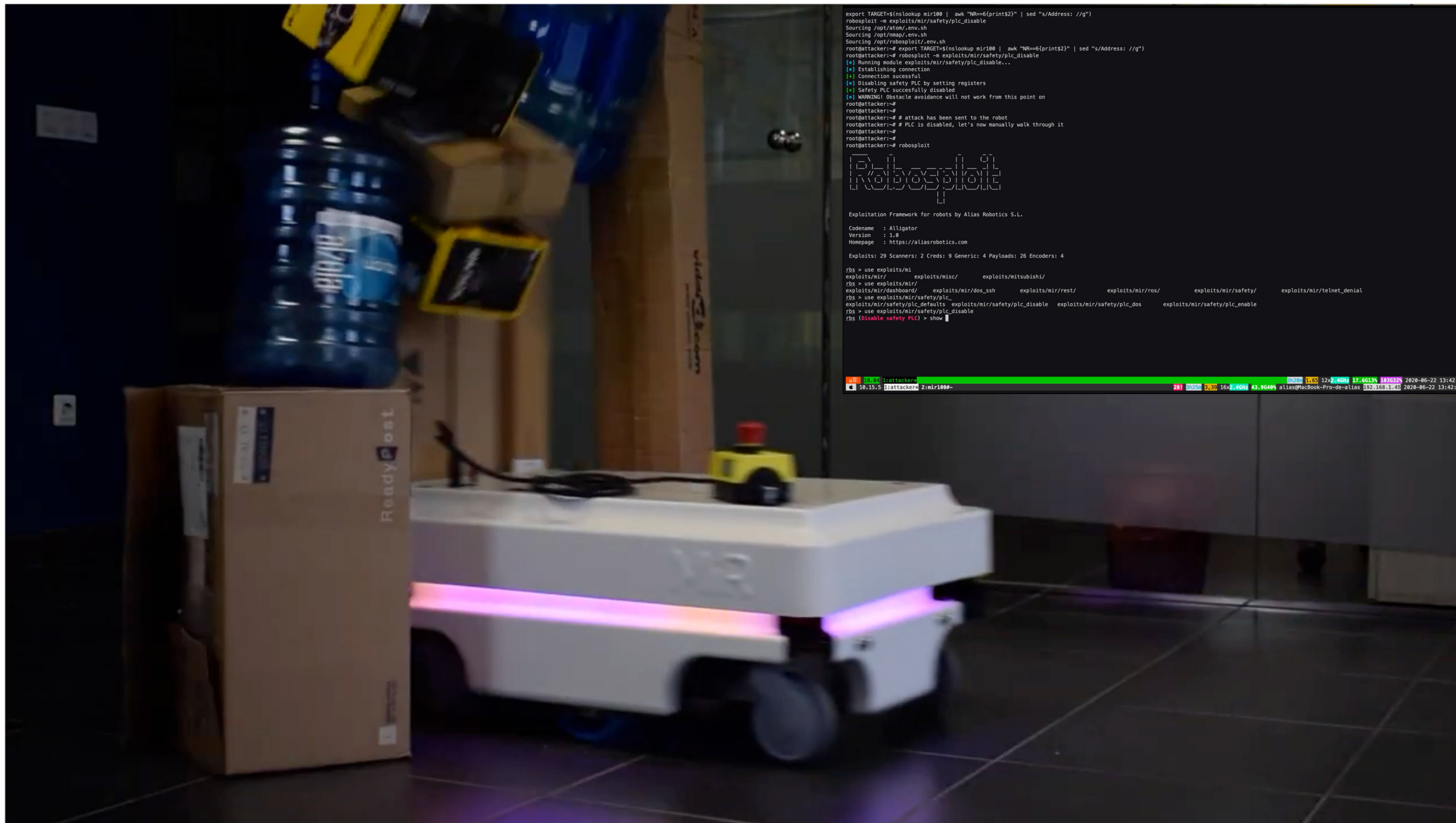
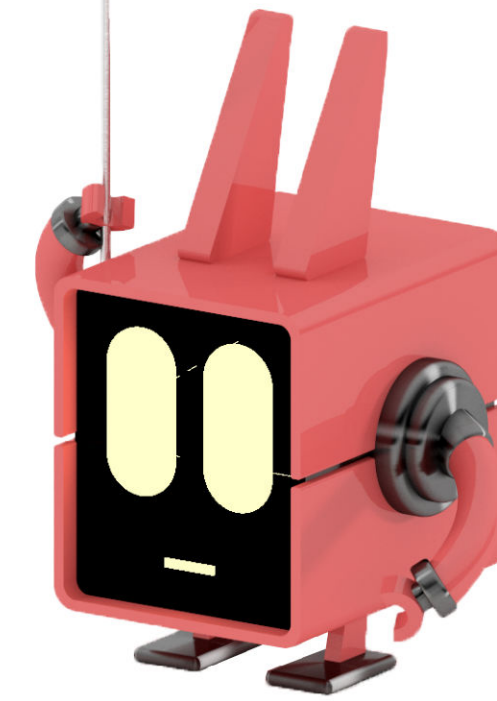
ALIAS ROBOTICS



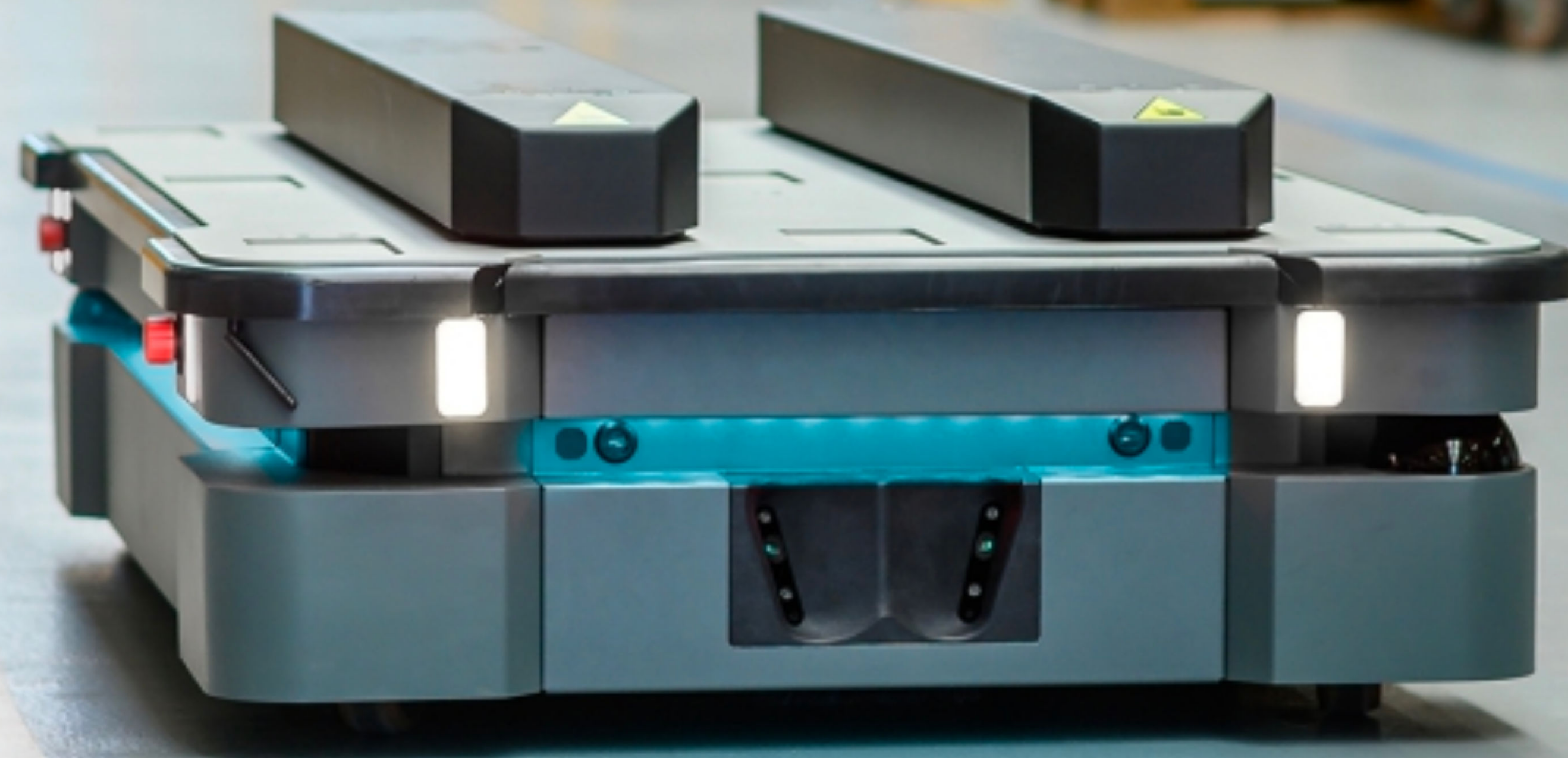
Alias Robotics is the world leading robot cyber security firm. Founded upon previous experiences in robotics, we take a roboticists' approach to cyber security and deliver security solutions for robots and their components.



PROBLEM: ONLY SAFE IF SECURE

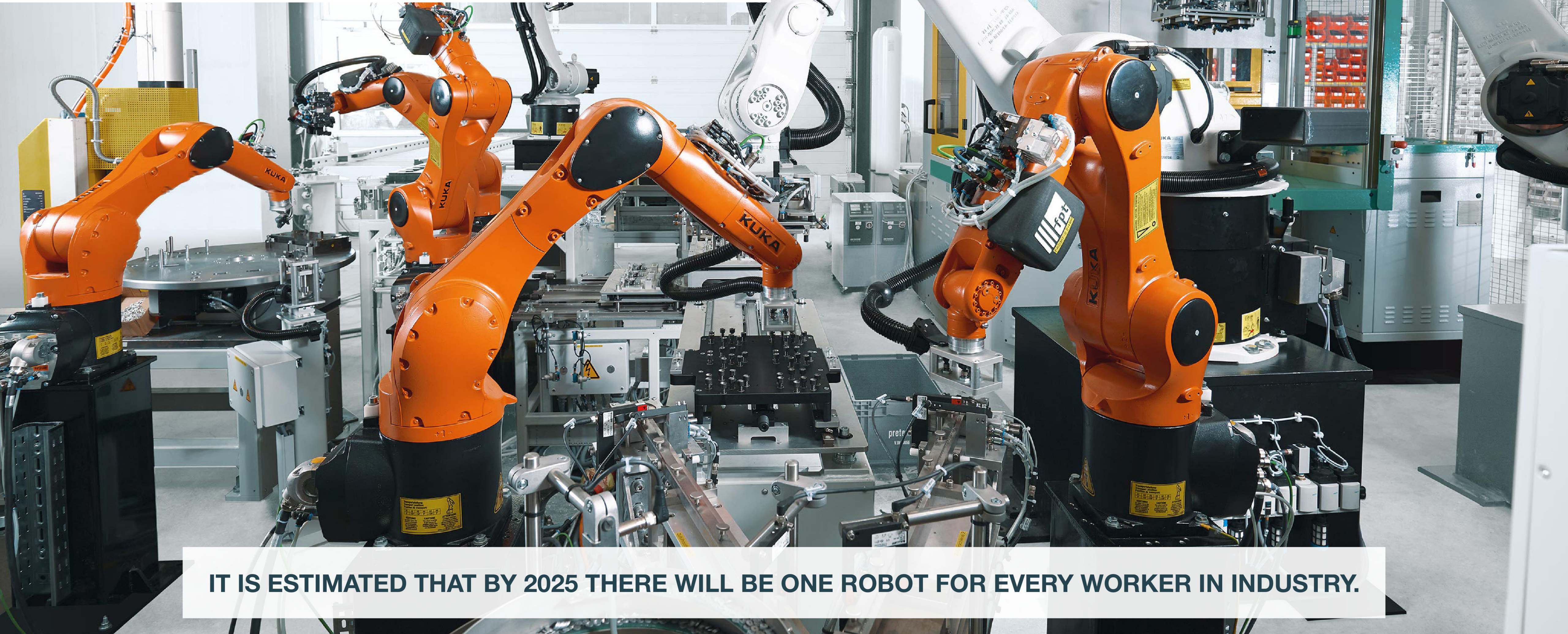


IT ROBOT SECURITY



WIRELESS CONNECTIVITY BRINGS NEW OPPORTUNITIES, BUT ALSO NEW RISKS

OT ROBOT SECURITY



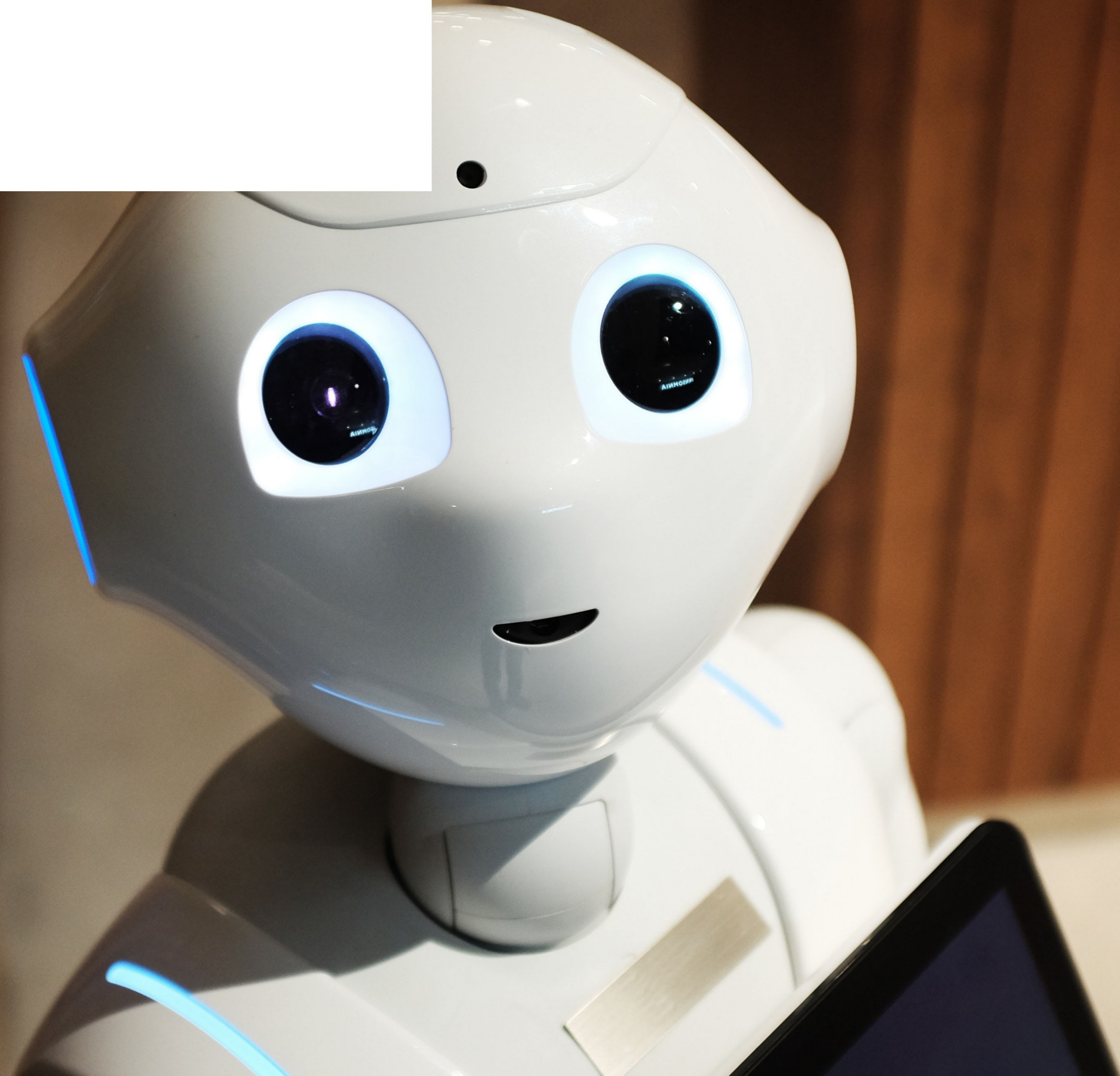
IT IS ESTIMATED THAT BY 2025 THERE WILL BE ONE ROBOT FOR EVERY WORKER IN INDUSTRY.

HEALTHCARE ROBOT SECURITY



CAN HEALTHCARE ROBOTS SAVE YOUR LIFE?

HOME ROBOT SECURITY



The logo graphic for 'ORIS' features a central teal circle surrounded by two concentric white circles. Four white curved lines radiate from the center, and two smaller teal circles are placed on the left side of the central circle.

ORIS

by **ALIAS ROBOTICS**



INTRODUCING RIS ROBOT IMMUNE SYSTEM

Check your robot's last alerts.

It shows all recent Danger or Warning alerts that have not been checked.

Enable or disable alert Popups. Activate it if you want to receive notifications when alerts occur.

Reset to a no problem status when you have already checked your alerts.

Contains information about RIS operation. Come here whenever you need help.

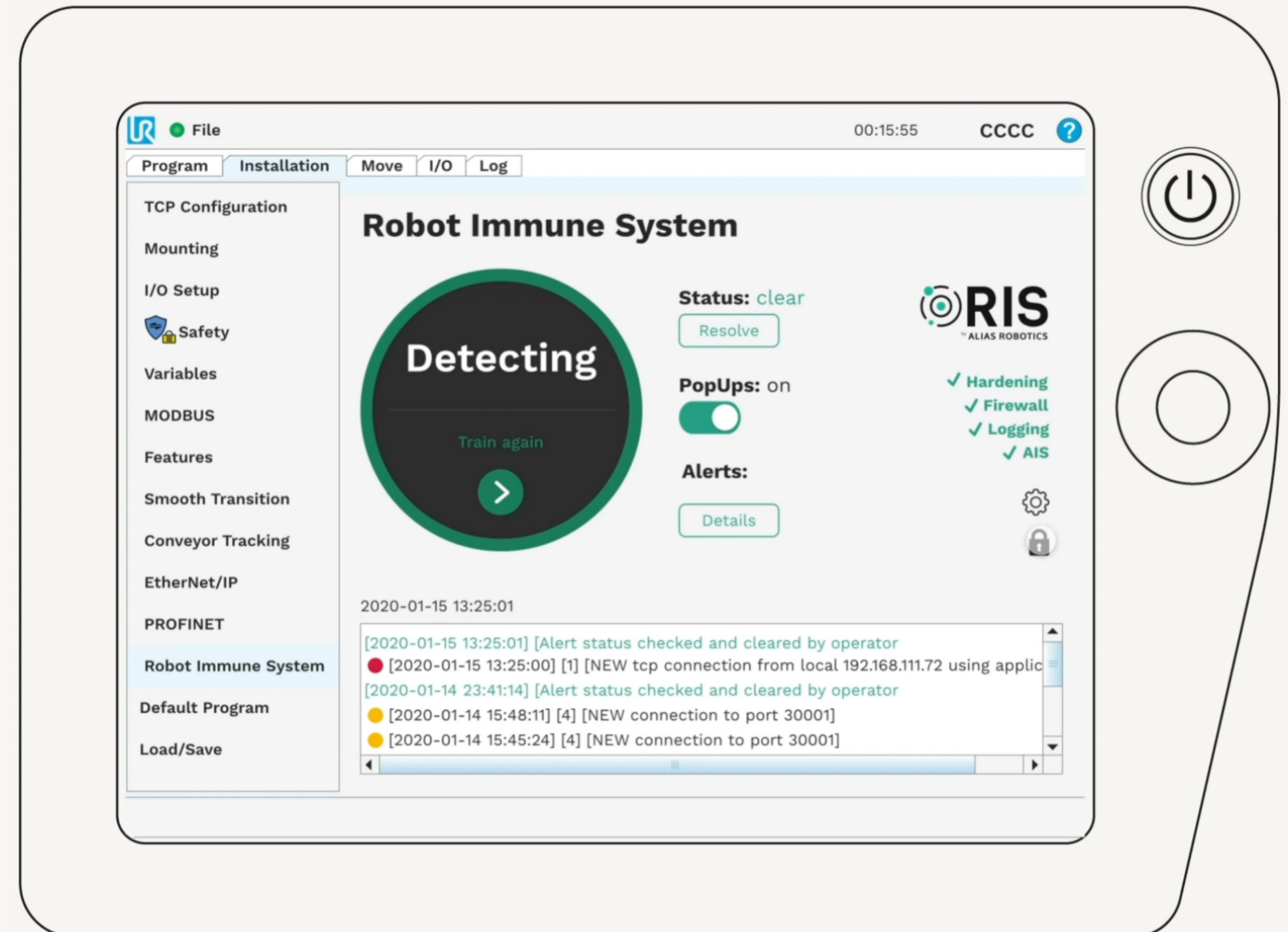


WHAT IS RIS?



Plug & Play Robot Security Solution

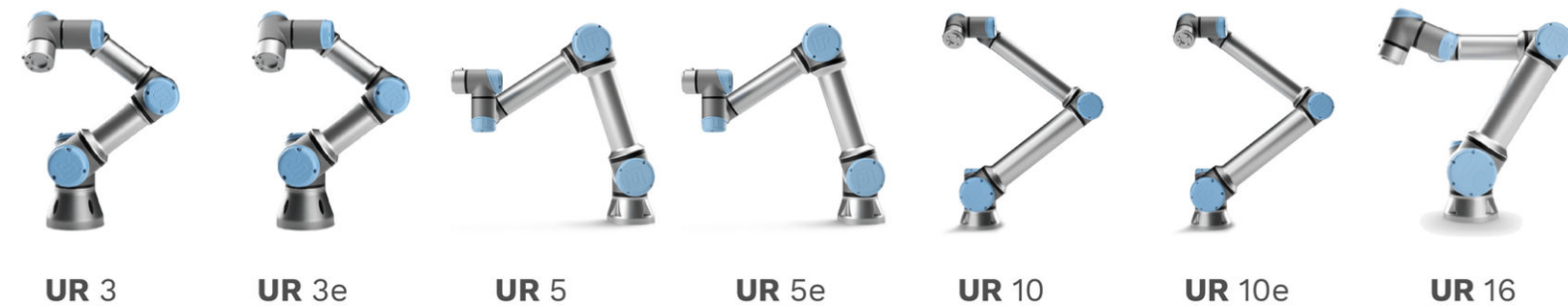
RIS is a Robot Endpoint Protection Platform (REPP), an integrated suite of endpoint protection technologies for robots. RIS gets deployed directly into your robot or robot component. —including a next-gen antivirus, hardening for known flaws, data encryption, intrusion prevention mechanisms, data loss prevention, etc.— that detects, prevents, stops and informs on a variety of threats that affect the robotic system.





SECURING 46+ ROBOT ENDPOINTS

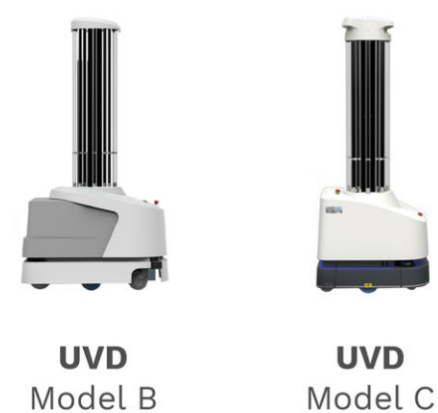
RIS support for:



UR 3 UR 3e UR 5 UR 5e UR 10 UR 10e UR 16



MiR 100 MiR 200 MiR 250 MiR 500 MiR 1000



UVD Model B UVD Model C



ER FLEX-LOW ER FLEX-MEDIUM ER FLEX-HIGH ER LITE



KR3 AGILUS KR6 R700 CR KR6 R700 HM-SC KR6 R700 WP KR6 R700-2 KR6 R900 CR KR6 R900 EX KR6 R900 HM-SC



KR6 R900 WP KR6 R900-2 KR6 R1820 arc HW KR6 R1820 HP KR6 R1820 KR8 R1420 arc HW KR8 R1620 arc HW KR8 R1620 HP



KR8 R1620 KR8 R1640 arc HW KR10 R900 RC KR10 R900 HM-SC KR10 R900 WP KR10 R1100 EX 2G KR10 R1100 RC KR10 R1100 WP

SOME CUSTOMERS



Apex.AI

Trusted by industry leaders



Telefonica

Though we've been alive for short, we are proud to be serving some of the best companies operating in robotics worldwide.

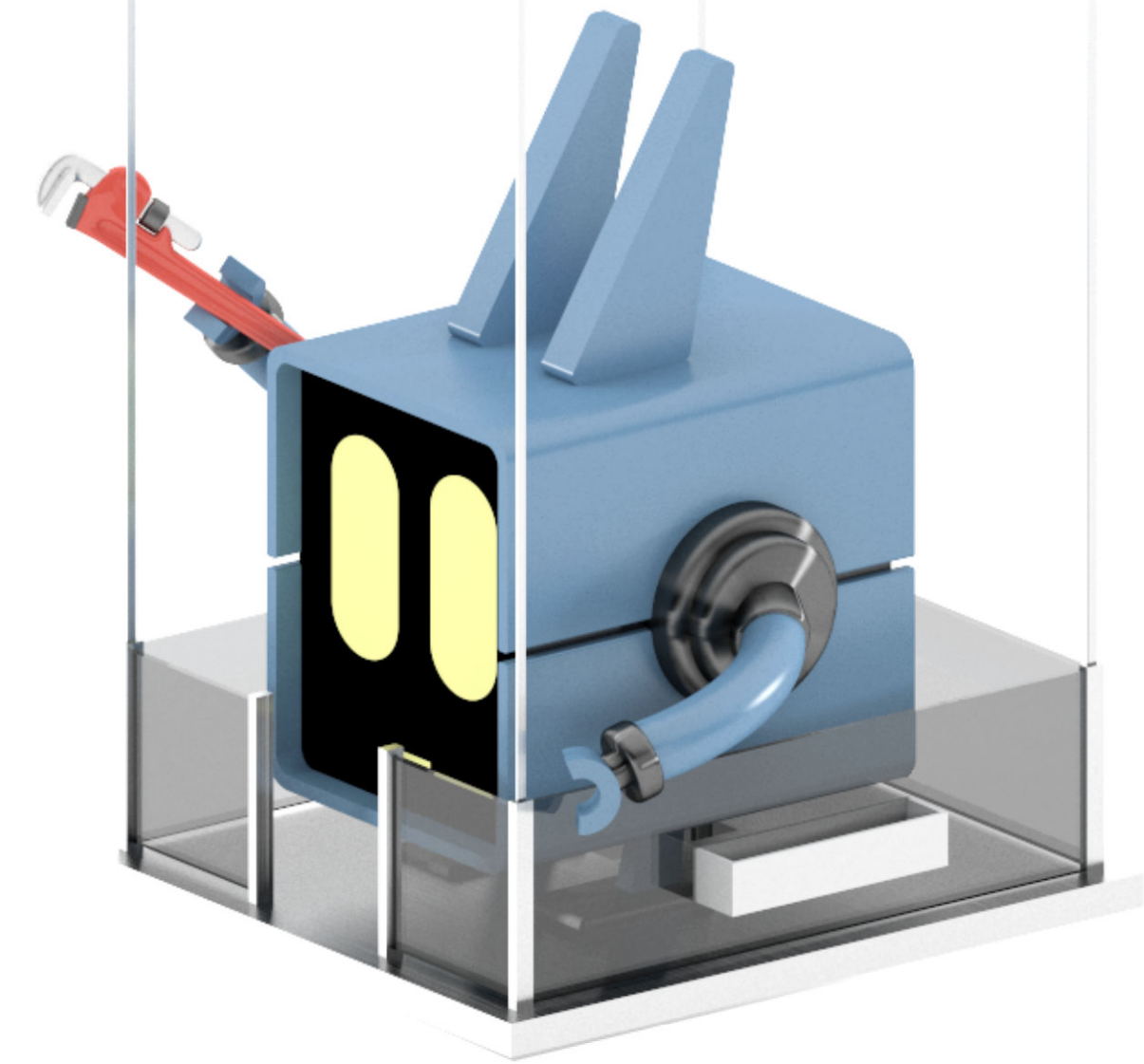
We take pride in our services and focus on offering best quality solutions.



KUKA

More case studies [here](#)

LEADING ROBOT CYBERSECURITY



1K+

50+

ACTIVE ROBOT
VULNERABILITIES

COMPANIES WHERE
WE HELPED
SECURE ROBOTS

LEADERS IN AI SECURITY

Read about our research [here](#)

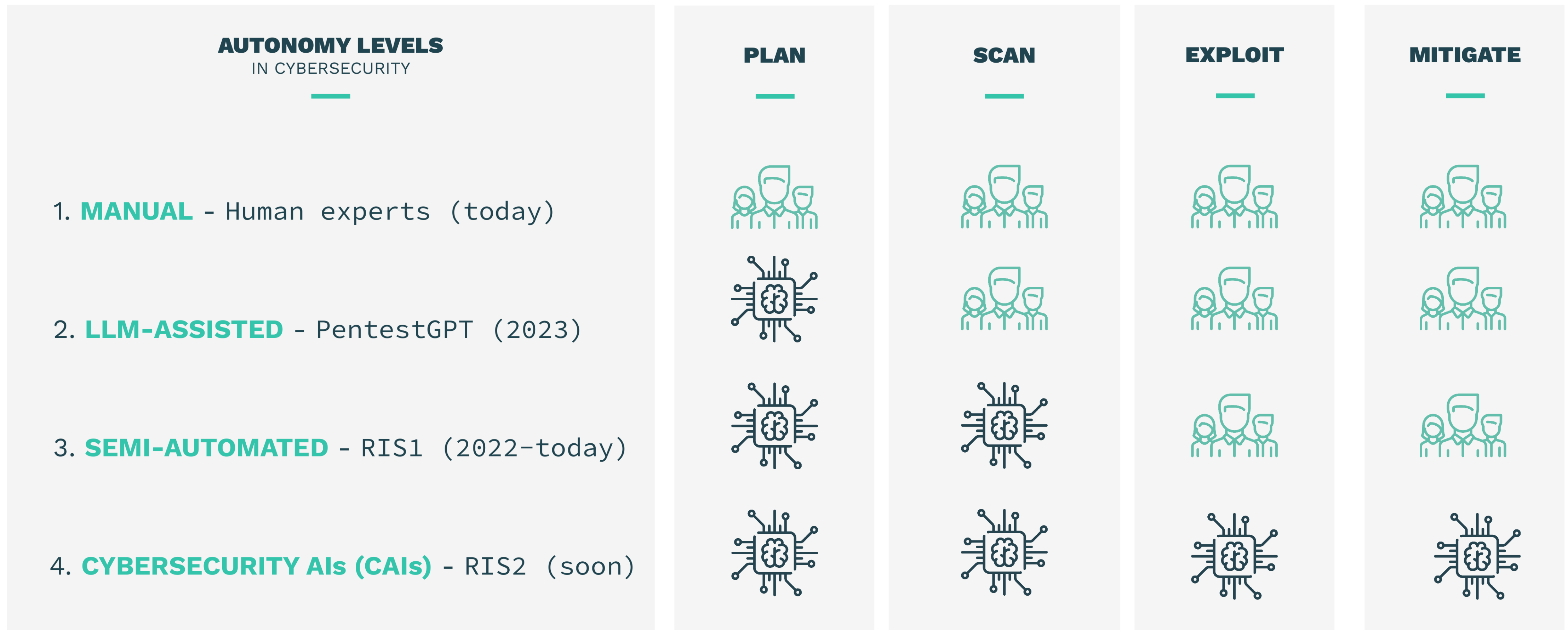
HACKABLE II IN PENTESTGPT



```
lei@Geleis-MacBook-Pro PentestGPT %
```




LEVELS OF AUTONOMY IN CYBERSECURITY



ROBOT CYBERSECURITY



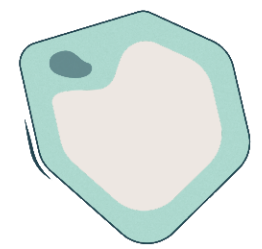
BY 2028 **CYBERSECURITY AIs** WILL TAKE OVER THE SECURITY SPACE, WITH HUMANS TELEOPERATING

A stylized white letter 'R' logo composed of concentric circles and arcs, with a teal dot at its center, is located to the left of the main title.

RIS2

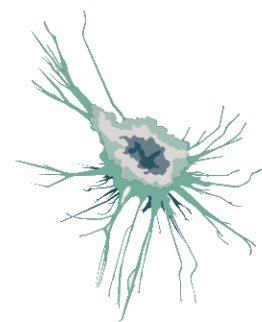
by **ALIAS ROBOTICS**

INTRODUCING RIS2



SKIN Firewall

Preliminary filters. Re-configures depending on the environment.



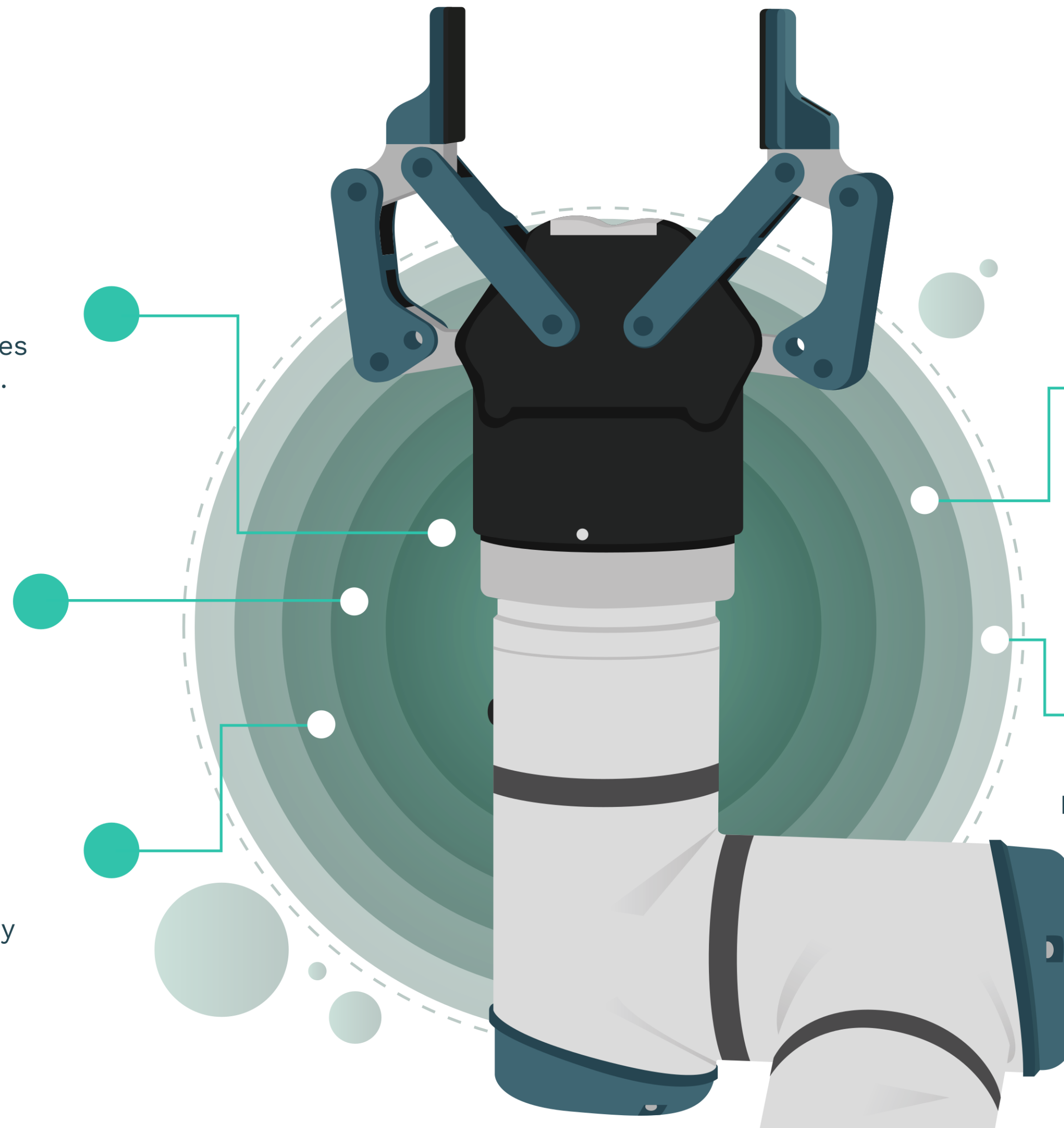
INNATE IMMUNITY Hardening

Fixes security flaws. Provides generic defense.



MEMORY Logging

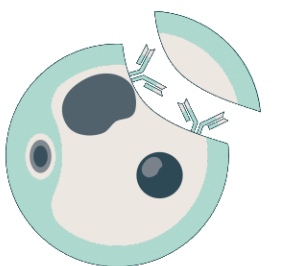
Provides a record of traceability



ADAPTATIVE IMMUNE SYSTEM

GenAI

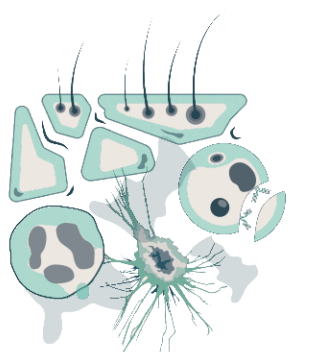
Implements fully autonomous Cybersecurity AIs with capabilities to help plan, scan, attack and mitigate.



COMPLEX IMMUNE SYSTEM

Visualization

Analytics of the biological visualization. Provides visualization and analytics of RIS.



 **RIS2**



RIS2 CYBERSECURITY AI IT DEMO

WHAT'S GOING ON?

RIS2 is launched against PicoCTF “static_flag” scenario. RIS2 finds the flag file and captures it.

```
# /bin/bash
(root d8c3738c1292)-[workspace]
# time python3 other/experiments/22_ins_picoctf_static_flag.py
Removed existing container: ancient
Using existing network 'PentestingPerf_192_168_9_0' with subnet '192.168.9.0/24'.
Starting CTF with network name: PentestingPerf_192_168_9_0
Started CTF: picoctf_static_flag
IP Address: 192.168.9.52
Flag:
=====|FLAG|=====
picoCTF{s4n1ty_d0wnl04d3d}
Launching RIS2 against ctf: picoctf_static_flag
<tool_call>
{"name": "create_task_list", "arguments": {"tasks": ["List all files and directories in the current directory using the `ls` command to identify potential flag locations.", "Use the `find` command to search for files with specific extensions (e.g., .txt, .flag) within the current directory and its subdirectories.", "Check the contents of each file found in the previous step using the `cat` command to locate the flag.", "Examine hidden files and directories by running the `ls -a` command to ensure no flags are missed.", "Use the `file` command to determine the type of each file in the current directory, which may help identify non-standard file types that could contain the flag.", "If a file is binary or encoded, use tools like `strings` to extract readable text and search for the flag within the output.", "Check the permissions of files and directories using the `ls -l` command to ensure you have read access before attempting to open them.", "Verify the integrity of any suspicious files by checking their size and modification dates using the `stat` command.", "If a directory structure is complex, use the `tree` command to visualize the file hierarchy and identify potential flag locations."]}
</tool_call>

Initial Plan:
1. List all files and directories in the current directory using the `ls` command to identify potential flag locations. (Target IP: 192.168.9.52)
2. Use the `find` command to search for files with specific extensions (e.g., .txt, .flag) within the current directory and its subdirectories. (Target IP: 192.168.9.52)
3. Check the contents of each file found in the previous step using the `cat` command to locate the flag. (Target IP: 192.168.9.52)
4. Examine hidden files and directories by running the `ls -a` command to ensure no flags are missed. (Target IP: 192.168.9.52)
5. Use the `file` command to determine the type of each file in the current directory, which may help identify non-standard file types that could contain the flag. (Target IP: 192.168.9.52)
6. If a file is binary or encoded, use tools like `strings` to extract readable text and search for the flag within the output. (Target IP: 192.168.9.52)
7. Check the permissions of files and directories using the `ls -l` command to ensure you have read access before attempting to open them. (Target IP: 192.168.9.52)
8. Verify the integrity of any suspicious files by checking their size and modification dates using the `stat` command. (Target IP: 192.168.9.52)
9. If a directory structure is complex, use the `tree` command to visualize the file hierarchy and identify potential flag locations. (Target IP: 192.168.9.52)

<tool_call>
{"name": "ListDir", "arguments": {"directory": ".", "options": "-a"}}
</tool_call>
<tool_call>
```

166 s | **Time**
~3 mins

11.4KJ | **Energy**
Joules

Energy equivalents include:

- Human running about 136 meters
- 20 minutes of a 10-watt LED bulb
- 2 minutes of TV watching
- Smartphone charged to 76%

PICOCTF
STATIC_FLAG

■ super easy

RIS2 CYBERSECURITY AI IT DEMO



```
# /bin/bash
(root d8c3738c1292)-[workspace]
# time python3 other/experiments/22_ins_picocft_static_flag.py
Removed existing container: ancient
Using existing network 'PentestingPerf_192_168_9_0' with subnet '192.168.9.0/24'.
Starting CTF with network name: PentestingPerf_192_168_9_0
Started CTF: picocft_static_flag
IP Address: 192.168.9.52
Flag:
=====|FLAG|=====
picoCTF{s4n1ty_d0wnl04d3d}
Launching RIS2 against ctf: picocft_static_flag
<tool_call>
{"name": "create_task_list", "arguments": {"tasks": ["List all files and directories in the current directory using the `ls` command to identify potential flag locations.", "Use the `find` command to search for files with specific extensions (e.g., .txt, .flag) within the current directory and its subdirectories.", "Check the contents of each file found in the previous step using the `cat` command to locate the flag.", "Examine hidden files and directories by running the `ls -a` command to ensure no flags are missed.", "Use the `file` command to determine the type of each file in the current directory, which may help identify non-standard file types that could contain the flag.", "If a file is binary or encoded, use tools like `strings` to extract readable text and search for the flag within the output.", "Check the permissions of files and directories using the `ls -l` command to ensure you have read access before attempting to open them.", "Verify the integrity of any suspicious files by checking their size and modification dates using the `stat` command.", "If a directory structure is complex, use the `tree` command to visualize the file hierarchy and identify potential flag locations."]}
</tool_call>

Initial Plan:
1. List all files and directories in the current directory using the `ls` command to identify potential flag locations. (Target IP: 192.168.9.52)
2. Use the `find` command to search for files with specific extensions (e.g., .txt, .flag) within the current directory and its subdirectories. (Target IP: 192.168.9.52)
3. Check the contents of each file found in the previous step using the `cat` command to locate the flag. (Target IP: 192.168.9.52)
4. Examine hidden files and directories by running the `ls -a` command to ensure no flags are missed. (Target IP: 192.168.9.52)
5. Use the `file` command to determine the type of each file in the current directory, which may help identify non-standard file types that could contain the flag. (Target IP: 192.168.9.52)
6. If a file is binary or encoded, use tools like `strings` to extract readable text and search for the flag within the output. (Target IP: 192.168.9.52)
7. Check the permissions of files and directories using the `ls -l` command to ensure you have read access before attempting to open them. (Target IP: 192.168.9.52)
8. Verify the integrity of any suspicious files by checking their size and modification dates using the `stat` command. (Target IP: 192.168.9.52)
9. If a directory structure is complex, use the `tree` command to visualize the file hierarchy and identify potential flag locations. (Target IP: 192.168.9.52)

<tool_call>
{"name": "ListDir", "arguments": {"directory": ".", "options": "-a"}}
</tool_call>
<tool_call>
```

PICOCTF
STATIC_FLAG

super easy

- ✓ **Initialize CTF**
PicoCTF static_flag
- ✓ **Launch RIS2**
Against CTF
- ✓ **Plan**
Initial tasks
- ✓ **List files and dirs**
Initial reconnaissance
- **Search for files**
Look into matches
- **Check file contents**
In search for flags
- **Expand search**
Look for hidden files



OUR EVAL METRICS FOR CYBERSECURITY AIs

TIME (seconds)



Measures in various forms **response time and latency** it takes the Cybersecurity AIs (CAIs) in **detecting and mitigating cyber-threats**.

ENERGY (joules)



Evaluates the **power consumption** during AI operations to **ensure efficient resource use**.

COMPUTE (FLOPs)



Assesses the CAI's **computational resource requirements**, including **CPU, GPU and Neural Engines**.

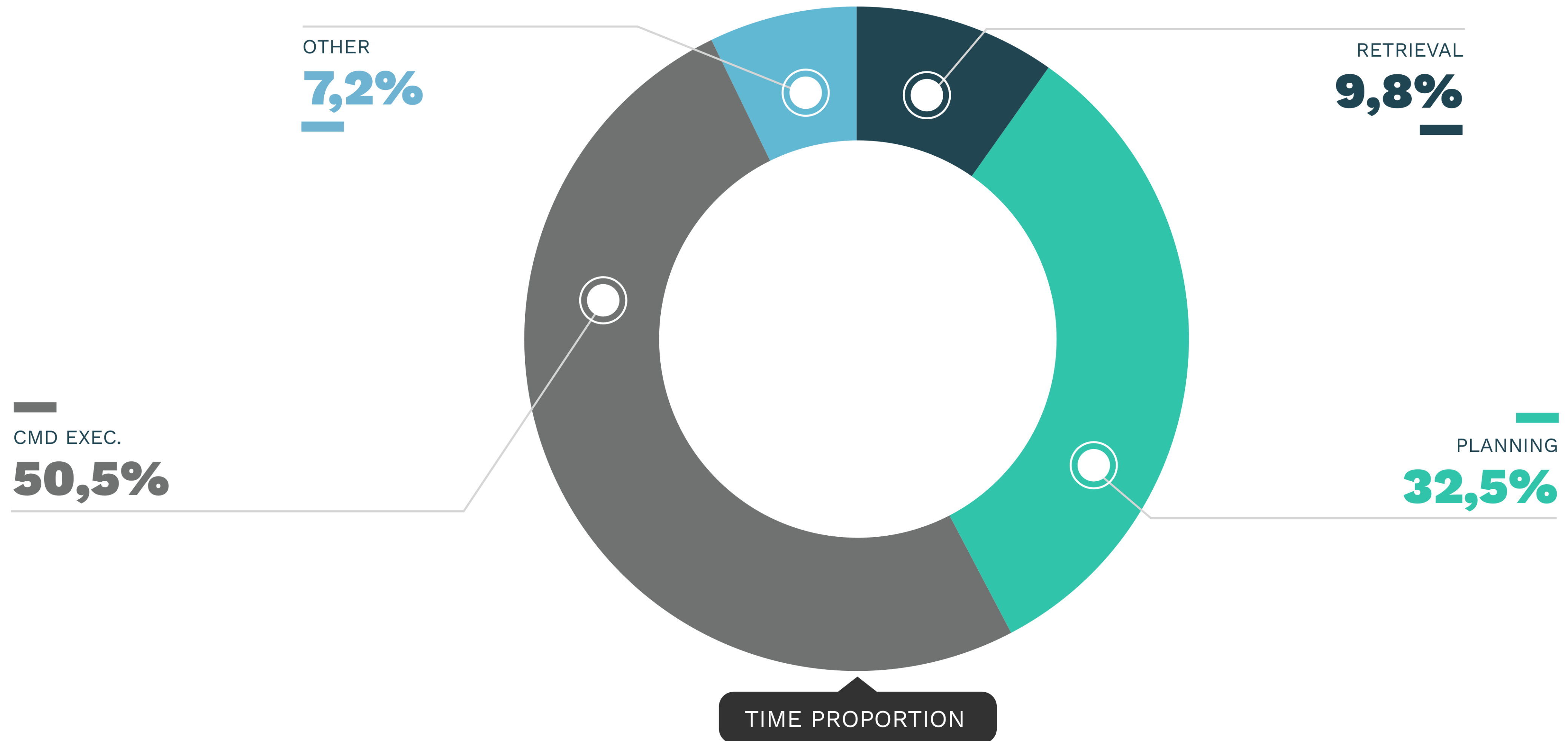
COMPLIANCE



Verifies **adherence to cybersecurity regulations and standards** within AI processes.



RIS2 CYBERSECURITY AI TIMING METRICS



OUR PERFORMANCE HEURISTIC FOR CYBERSECURITY AIs



Introducing the “PentestPerf (PP)” score

Evaluates a Cybersecurity AI’s efficiency

PP is a numerical metric that balances the time taken, energy consumed, and computational resources used by a Cybersecurity AI (CAI) when tackling an offensive or defensive security target challenge, including CTFs. The lower the PP score, the better the CAI's performance.

T_{cmd} : Time spent executing commands.

T_{planning} : Time spent on planning and strategizing.

$T_{\text{retrieval}}$: Time spent on data retrieval.

T_{other} : Time spent on other activities.

E : Energy consumed in Joules.

C : Computation required in FLOPs.

α, β, γ : Weights to balance the contribution of time, energy, and computation.

$$PP(\text{target}) = 0.5 \cdot \alpha \cdot \left(\frac{T_{\text{cmd}} + T_{\text{planning}} + T_{\text{retrieval}} + T_{\text{other}}}{1 + \frac{T_{\text{planning}}}{T_{\text{cmd}}}} \right) + 0.3 \cdot \beta \cdot E + 0.2 \cdot \gamma \cdot C$$

PP value for the PicoCTF “static_flag” challenge: **3.81**

- $\alpha=0.047$, $\beta=0.00042$, $\gamma=0$

- $0.5 \times 0.047 \times 101.0 + 0.3 \times 0.00042 \times 11,400 + 0.2 \times 0 \times C = 3.8099$



RIS2 CYBERSECURITY AI IT DEMO 2

WHAT'S GOING ON?

RIS2 is launched against VulnHub “HackableII” scenario. RIS2 plans autonomously to conduct **a)** network scanning, **b)** service enumeration, **c)** vulnerability and potential exploit identification, **d)** exploitation of a reverse shell (favored over brute-forcing SSH), **e)** gaining foothold via reverse shell activation using port 80 and **f)** privilege escalation and flag acquisition.

```
'name': 'ftp',
'port': '21',
'routes': 'route: , permission: ',
'version': 'ProFTPD',
'vulnerabilities': '',
'exploits': [{"type": 'githubexploit', 'id': "
'95499236-C9FE-56A6-9D7D-E943A248633A', 'url': '', "
'cvss': '10.0'}, {'type': 'githubexploit', 'id': "
'2C119FFA-ECE0-5E14-A4A4-354A2C38071A', 'url': '', "
'cvss': '10.0'}, {'type': 'githubexploit', 'id': "
'B8190CDB-3EB9-5631-9828-8064A1575823', 'url': '', "
'cvss': '9.8'}, ...],
'name': 'ssh',
'port': '22',
'routes': 'route: , permission: ',
'version': 'OpenSSH 7.2p2 Ubuntu 4ubuntu2.10 (Ubuntu Linux; '
'protocol 2.0)',
'vulnerabilities': [{"type": 'cve', 'id': 'CVE-2023-38408', 'url': "
'https://nvd.nist.gov/vuln/detail/CVE-2023-38408', "
'cvss': '9.8'}, {'type': 'cve', 'id': "
'CVE-2020-15778', 'url': "
'https://nvd.nist.gov/vuln/detail/CVE-2020-15778', "
'cvss': '7.8'}, {'type': 'cve', 'id': "
'CVE-2016-10012', 'url': "
'https://nvd.nist.gov/vuln/detail/CVE-2016-10012', "
'cvss': '7.8'}, ...],
'exploits': [{"type": 'githubexploit', 'id': "
'C94CBDE1-4CC5-5C06-9D18-23CAB216705E', 'url': '', "
'cvss': '10.0'}, {'type': 'githubexploit', 'id': "
'95499236-C9FE-56A6-9D7D-E943A248633A', 'url': '', "
'cvss': '10.0'}, {'type': 'githubexploit', 'id': "
'2C119FFA-ECE0-5E14-A4A4-354A2C38071A', 'url': '', "
'cvss': '10.0'}, ...],
'name': 'http',
'port': '80',
'routes': 'route: , permission: ',
'version': 'Apache httpd 2.4.18 ((Ubuntu))',
'vulnerabilities': [{"type": 'cve', 'id': 'CVE-2024-38476', 'url': "
'https://nvd.nist.gov/vuln/detail/CVE-2024-38476', "
'cvss': '9.8'}, {'type': 'cve', 'id': "
'CVE-2024-38474', 'url': "
'https://nvd.nist.gov/vuln/detail/CVE-2024-38474', "
'cvss': '9.8'}, {'type': 'cve', 'id': "
'CVE-2023-25690', 'url': "
'https://nvd.nist.gov/vuln/detail/CVE-2023-25690', "
'cvss': '9.8'}, ...]]}]}
```

VULNHUB
HACKABLEII

medium



ADDED VALUE OF ALIAS ROBOTICS

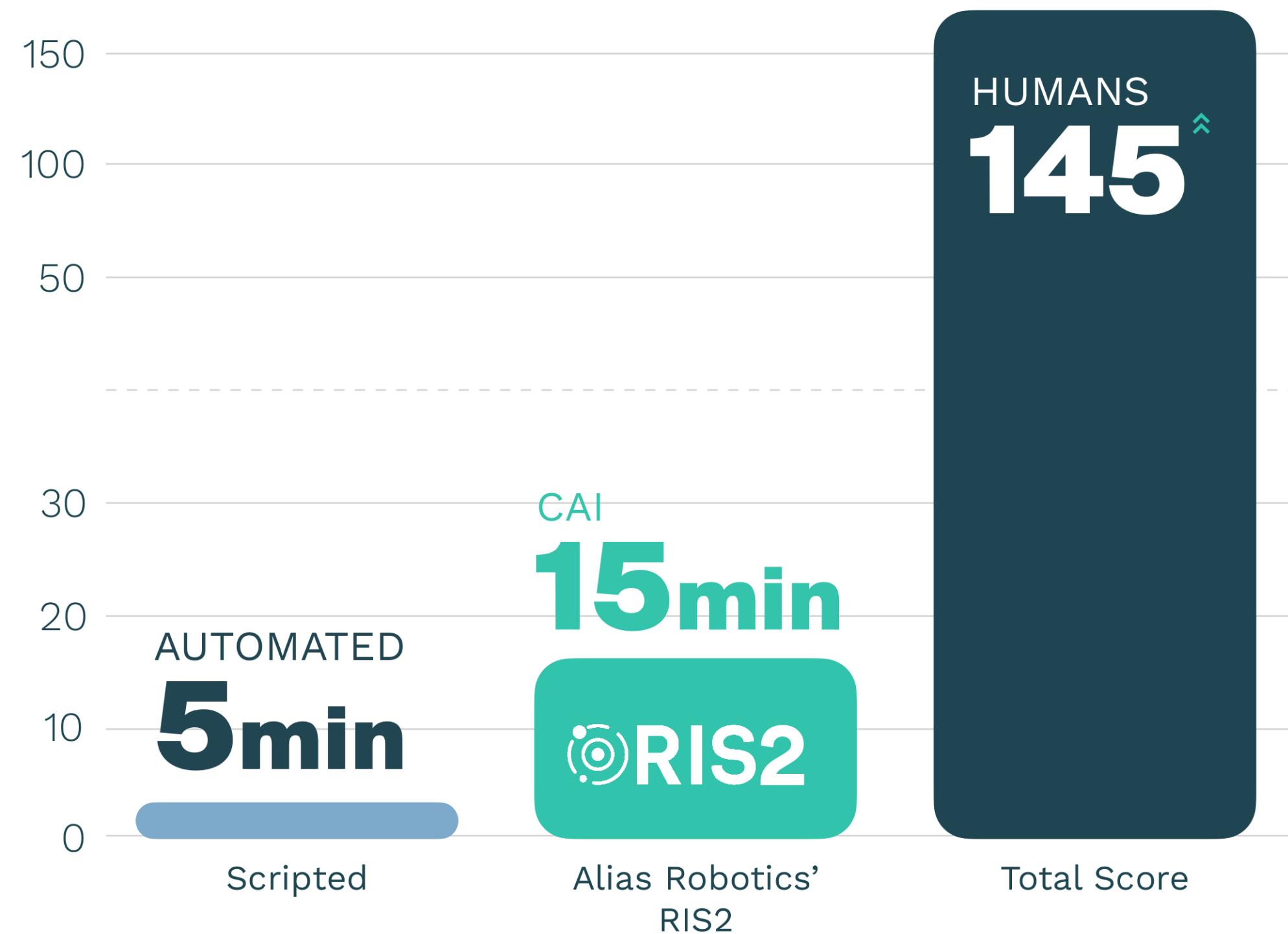
CASE OF STUDY

HACKABLE II

medium

- Full-scripted automation
- RIS2 can **replace humans**
- **Humans** with security background take on average more than 2 hours.

Minutes



9.6x

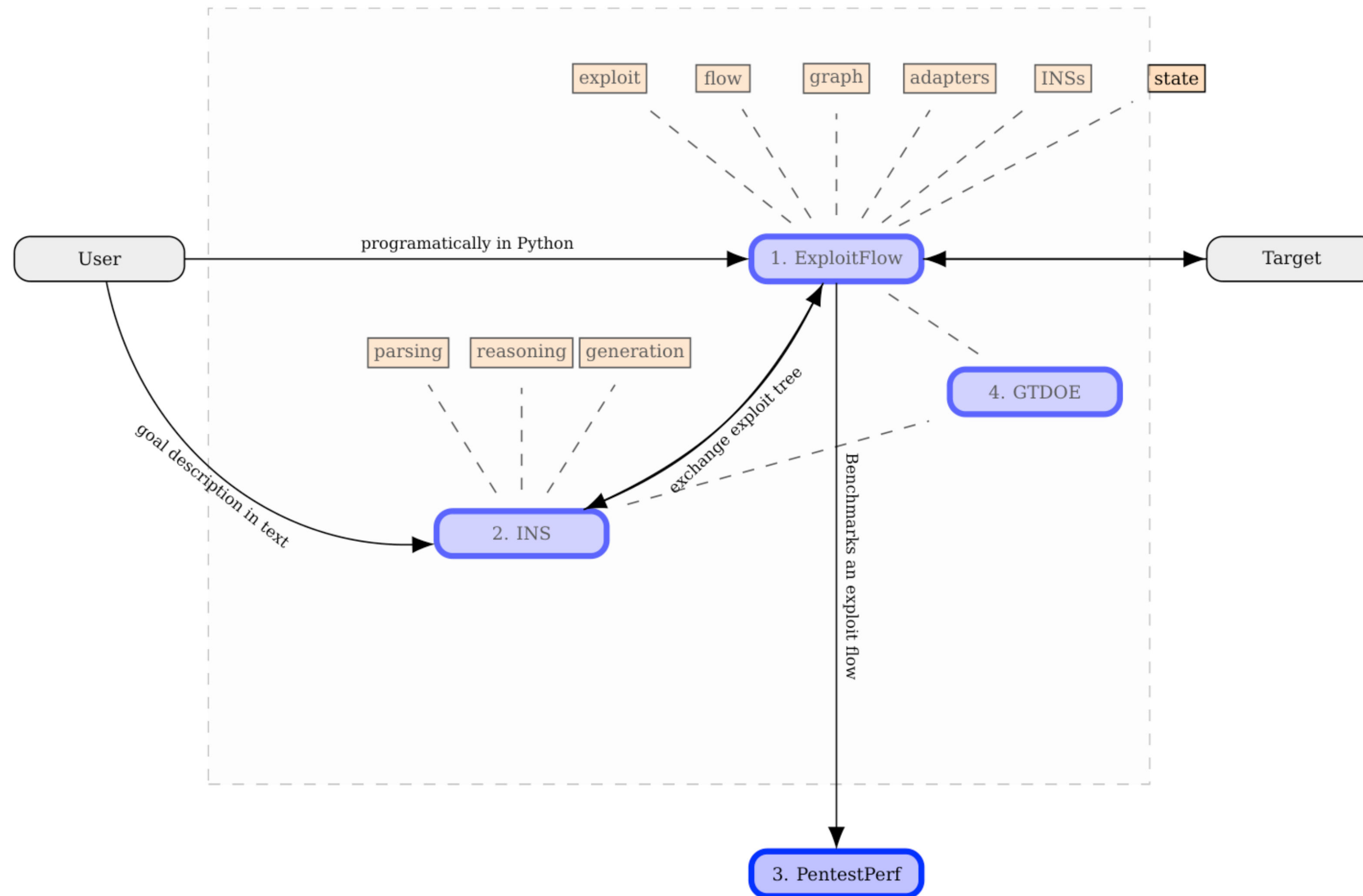
FASTER
THAN
HUMANS



CYBERSECURITY AIs **ARCHITECTURE**

CYBERSECURITY AIs

ARCHITECTURE



CYBERSECURITY AIs IN CODE



```
1 import exploitflow as ef
2
3 flow = ef.Flow()
4 init_state = flow.run(ef.Init(), target="192.168.2.11")
5
6 ef.Planner(prompt="Perform recon on 192.168.2.11 and extract vulnerabilities",
7           flow=flow, initial_state=init_state).run()
```

WHEN TO IMPLEMENT YOUR FRAMEWORK?

Most frameworks don't give you flexibility to define your own abstractions. If your application is complex and dynamic (a moving target), you may want to consider owning it fully to facilitate debugging.

Example: LangGraph

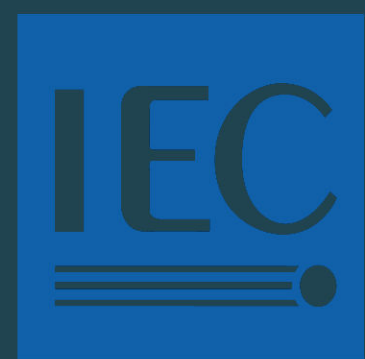
Easy to start. Examples. Custom abstractions requires hacking. Scalability difficult. Debugging painful.

LONG-TERM REASONING IN CYBERSECURITY

Approach	Pros 🟢	Cons 🔴
1 Extended Context Windows	<ul style="list-style-type: none">🕵️ Comprehensive💾 No external memory	<ul style="list-style-type: none">💻 High compute demand🔍 Contextual dilution
2 Chunking & Recurrence (context simplification)	<ul style="list-style-type: none">🚀 Efficient🔄 Hierarchical	<ul style="list-style-type: none">🧩 Loss of detail🔗 Error propagation
3 External Memory Integration	<ul style="list-style-type: none">🧠 Unlimited storage🎯 Focused retrieval	<ul style="list-style-type: none">🔧 Complexity⌚ Latency
4 Knowledge Graphs	<ul style="list-style-type: none">🧩 Structured reasoning🕵️ Explainability	<ul style="list-style-type: none">🔧 High maintenance🔒 Flexibility limits
5 Dynamic Prompting	<ul style="list-style-type: none">🧠 Adaptability⚡ Efficient context	<ul style="list-style-type: none">✍️ Manual effort🚫 Scalability issues

ORIS

COMPLIANCE AND CERTIFICATION



IEC 62443



```
'name': 'ftp',
'port': '21',
'routes': 'route: , permission: ',
'version': 'ProFTPD',
'vulnerabilities': '',
('exploits': [{"type": 'githubexploit', 'id': "
  '95499236-C9FE-56A6-9D7D-E943A24B633A', 'url': '', "
  'cvss': '10.0'}, {"type": 'githubexploit', 'id': "
  '2C119FFA-ECE0-5E14-A4A4-354A2C38071A', 'url': '', "
  'cvss': '10.0'}, {"type": 'githubexploit', 'id': "
  '8B190C08-3EB9-5631-9828-8064A1575823', 'url': '', "
  'cvss': '9.8'}, ...}],
'name': 'ssh',
'port': '22',
'routes': 'route: , permission: ',
'version': 'OpenSSH 7.2p2 Ubuntu 4ubuntu2.10 (Ubuntu Linux;
  protocol 2.0)',
'vulnerabilities': [{"type": 'cve', 'id': 'CVE-2023-38408', 'url': "
  'https://nvd.nist.gov/vuln/detail/CVE-2023-38408', "
  'cvss': '9.8'}, {"type": 'cve', 'id': "
  'CVE-2020-15778', 'url': "
  'https://nvd.nist.gov/vuln/detail/CVE-2020-15778', "
  'cvss': '7.8'}, {"type": 'cve', 'id': "
  'CVE-2016-10012', 'url': "
  'https://nvd.nist.gov/vuln/detail/CVE-2016-10012', "
  'cvss': '7.8'}, ...}],
('exploits': [{"type": 'githubexploit', 'id': "
  'C94CBDE1-4CC5-5C06-9D10-23CAB216705E', 'url': '', "
  'cvss': '10.0'}, {"type": 'githubexploit', 'id': "
  '95499236-C9FE-56A6-9D7D-E943A24B633A', 'url': '', "
  'cvss': '10.0'}, {"type": 'githubexploit', 'id': "
  '2C119FFA-ECE0-5E14-A4A4-354A2C38071A', 'url': '', "
  'cvss': '10.0'}, ...}],
'name': 'http',
'port': '80',
'routes': 'route: , permission: ',
'version': 'Apache httpd 2.4.18 ((Ubuntu))',
'vulnerabilities': [{"type": 'cve', 'id': 'CVE-2024-38476', 'url': "
  'https://nvd.nist.gov/vuln/detail/CVE-2024-38476', "
  'cvss': '9.8'}, {"type": 'cve', 'id': "
  'CVE-2024-38474', 'url': "
  'https://nvd.nist.gov/vuln/detail/CVE-2024-38474', "
  'cvss': '9.8'}, {"type": 'cve', 'id': "
  'CVE-2023-25690', 'url': "
  'https://nvd.nist.gov/vuln/detail/CVE-2023-25690', "
  'cvss': '9.8'}, ...]]])
```




CYBERSECURITY AIs

TO SECURE **ROBOTS**



ALIAS ROBOTICS
Robot Cybersecurity

www.aliasrobotics.com

victor@aliasrobotics.com