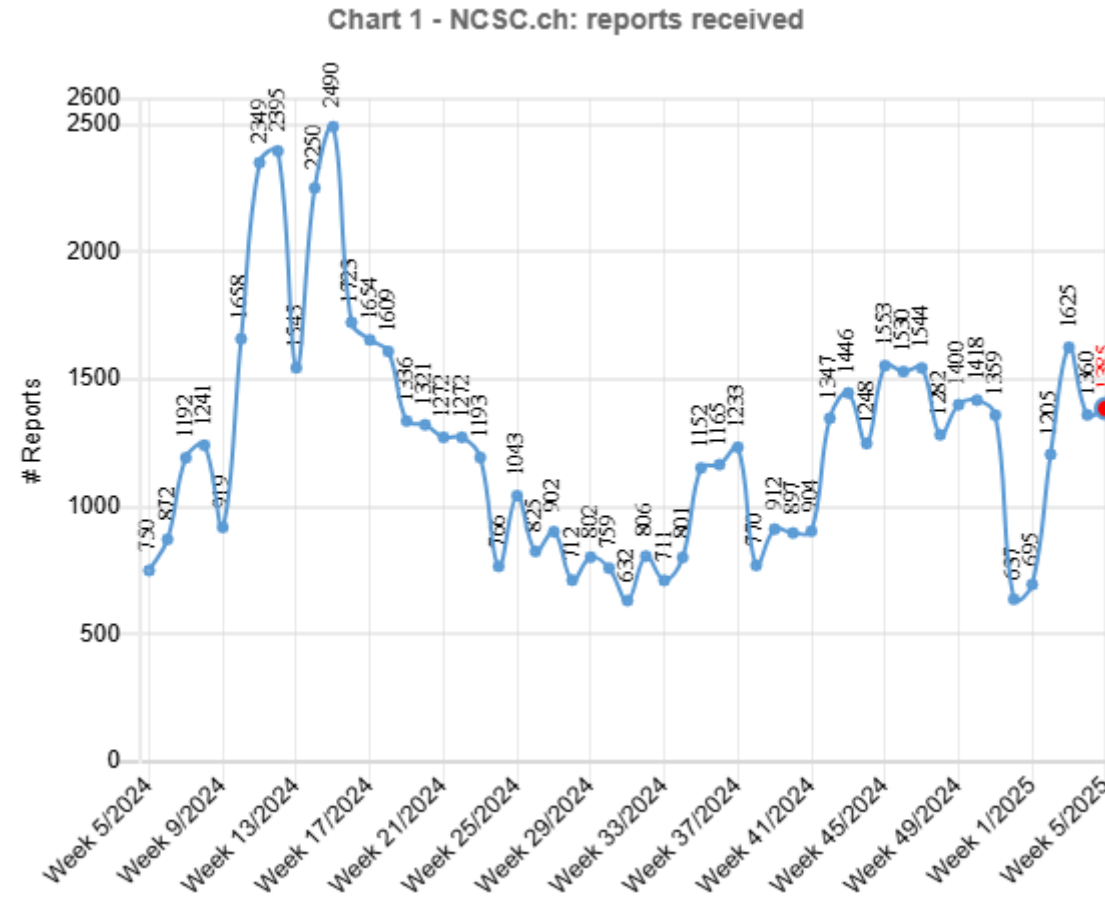# Towards a fully automated Blue Team at Locked Shields

Roland Meier
AMLD 2025, February 13, 2025

# Every week, the NCSC receives >1`000 reports about cyber incidents



Chart 1 - NCSC.ch: reports received

# Every week, the NCSC receives >1`000 reports about cyber incidents

### E-mails with malware in the name of debt collection agencies and health insurance companies

02.12.2024 - The NCSC is currently receiving numerous reports of e-mails that claim to come from a debt collection agency or a health insurance company. They concern an alleged claim or reminder. Do not click on the link, as this is an attempt to distribute malware to Windows users.

### Update: Even after the conclusion of the high-level conference on peace in Ukraine, the overload attacks on websites of organisations involved continue

17.06.2024 - As expected, the overload attacks continue even after the conclusion of the high-level conference on peace in Ukraine. The websites of the organisations involved in the conference are still being targeted. The National Cyber Security Centre is monitoring the situation and is in contact with the organisations concerned.
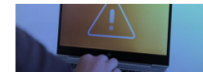
### Critical vulnerability in Palo Alto firewalls

18.04.2024 - The NCSC warns of the security vulnerability in Palo Alto's Next-Generation Firewall (NGFW). These firewalls are mainly used by companies and public authorities. They have a critical vulnerability that is already being exploited by cyber criminals. The attackers exploit the vulnerability to execute commands. The NCSC has already received corresponding reports from organisations in Switzerland. The NCSC recommends installing the security updates as quickly as possible or even reinstalling the NGFW if possible.

### Critical vulnerability in file transfer software «MOVEit»: Apply Patch quickly

02.06.2023 - The file transfer software called «MOVEit», which is mainly used by businesses, has a critical vulnerability that is already being exploited by cybercriminals. The attackers are exploiting the vulnerability to steal files from the file transfer software. The NCSC started to receive corresponding reports from organisations in Switzerland on 1 June. The NCSC recommends applying the security patch as quickly as possible.

### Update: Still over 2,000 unsecured Microsoft Exchange servers in Switzerland

01.12.2022 - Just over a fortnight ago, the NCSC called for the security patches provided by Microsoft to be installed in order to fix the ProxyNotShell vulnerability. Despite the urgency, there are still some operators that have failed to heed this call. Therefore, the NCSC has sent more than 2,000 registered letters to those ...ed, urging them to act now.

# There are not enough experts in this field

# AI is everywhere…



GENIUS X

THE
REVOLUTIONARY
GENIUS X WITH
ARTIFICIAL
INTELLIGENCE

Recognizes your brushing style. Guides you to brush
better every day.

SHOP NOW

… why not in cyber defense?

# THE IRISH TIMES

Mon, May 20, 2019

Dublin 14°C

## Cyber defenders fight hackers in high-tech Estonia war

Attacks on vital systems and fake news

Fri, Apr 28, 2017, 01:00

**Daniel McLaughlin** in Tallinn

💬 0

[Kaspersky]

The attack on the airbase began with a salvo of fake news. "A report appeared saying drones were using nerve gas," said Lauri Luht, crisis management chief for the cyber security department of Estonia's information system authority.

VBS / armasuisse / W+

# Cyber defenders fight hackers in high-tech Estonia war games

Attacks on vital systems and fake news are all part of Locked Shields exercise

🕐 Fri, Apr 28, 2017, 01:00

**Daniel McLaughlin** in Tallinn

💬 0



Locked Shields, now taking place in Estonia involving 20 teams from Europe and the US, is the world's most advanced live-fire cyber defence exercise. Photograph: Daniel McLaughlin

The attack on the airbase began with a salvo of fake news. "A report appeared saying drones were using nerve gas," said Lauri Luht, crisis management chief for the cyber security department of Estonia's information system authority.

VBS / armasuisse / W+T

**Locked Shields is the largest live-fire
global cyber defense exercise**

Dr. Roland Meier                                                                 9

Picture: NATO CCDCOE

# Locked Shields is the largest live-fire global cyber defense exercise



Picture: NATO CCDCOE

- Red Team vs. Blue Team exercise

  Attackers    Defenders
  1 Team    1 Team / country

- \> 1'000 experts from 30 nations

- \> 4'000 systems

- \> 2'500 attacks

# Locked Shields is organized by the NATO CCDCoE





Sponsoring Nations

Contributing Participants

# Each Blue Team is responsible for its network ("Gamenet")



e.g., a military air base

# The Gamenet consists of a large variety of systems

# Each Blue Team is responsible for its network ("Gamenet")



e.g., a military air base

3 main tasks:

- Perform initial hardening

- Defend against attacks

- Communicate with other teams

# Besides defending its network, a Blue Team needs to communicate with other teams

User Simulation Team      Read, address and respond to support tickets

Yellow Team               Provide periodic reports

White team                Voice or video calls through the Gamenet

Green Team                Gamenet status, reverting of devices

# The number of people required
# in a Blue Team continuously increases



| | 2010* | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Max. Number of Persons on a Blue Team | 10 | | 10 | 10 | 16 | 16 | 16 | 56 | 69 | 104 | | |
| Blue Teams | 6 | | 9 | 10 | 12 | 15 | 20 | 20 | 22 | 23 | | 22 |

■ Max. Number of Persons on a Blue Team    ■ Blue Teams

[Smeets, Max. "The role of military cyber exercises: A case study of Locked Shields." *CyCon 2022*]

# Research goals

How can automation / AI help for cyber defense?

*And eventually…*

What would it take to have a fully automated Blue Team

in a future iteration of Locked Shields?

# The history

# Four stages in the exercise

| Initial hardening | Monitoring & response | Reporting | Recovery |

**Prepare systems
before attacks start**

**Monitor systems while
they are under attack
and mitigate attacks**

**Respond to user requests
and provide incident reports**

**Restore and fix
compromised systems**

# Automated Blue Team framework overview

| Input | Control | Output |
|:---:|:---:|:---:|

# Automated Blue Team framework overview

Sensors (examples)

Network traffic → Sensor Input Interface

Event logs →

External monitoring →

Support tickets →

Device credentials →

Control

Output

# Automated Blue Team framework overview

Sensors (examples)

Actuators (examples)

Network traffic

Event logs

External monitoring

Support tickets

Device credentials

Sensor Input Interface

Control

Actuator Output Interface

Remote management

Firewall and routing rules

Device reset

Chat message

Report

# Automated Blue Team framework overview



Sensors (examples)

- Network traffic
- Event logs
- External monitoring
- Support tickets
- Device credentials

Sensor Input Interface

Situational Awareness

AI Engine

Control Logic

Actuator Output Interface

Actuators (examples)

- Remote management
- Firewall and routing rules
- Device reset
- Chat message
- Report

# Automated Blue Team framework overview

# Five tasks for AI

Identification / classification        What is it?

Categorisation                         What belongs together?

Assessment                             What is important?

Recommendation                         What should be done?

Prediction                             What will happen?

# AI for initial hardening

| Initial hardening | Monitoring & response | Reporting | Recovery |

Identification / classification

Categorisation

Find groups of similar devices

Assessment

Recommendation

Prediction

# AI for monitoring and response

| Initial hardening | Monitoring & response | Reporting | Recovery |

Identification / classification

Categorisation

Assessment

Recommendation

Prediction

Detect malicious network traffic

Detect malicious patterns in log files

# AI for reporting

| Initial hardening | Monitoring & response | Reporting | Recovery |
| --- | --- | --- | --- |

Identification / classification

Categorisation                          Link support tickets and monitoring alerts

Assessment                              Prioritise support tickets

Recommendation                          Formulate response to support tickets

Prediction                              Predict impact on the scoring

# AI for recovery

| Initial hardening | Monitoring & response | Reporting | Recovery |

Identification / classification      Find devices that need to be recovered

Categorisation      Find similar devices as a template for recovery

Assessment

Recommendation

Prediction

# Automated Blue Team framework overview



Sensors (examples)

Network traffic

Event logs

External monitoring

Support tickets

Device credentials

Sensor Input Interface

Situational Awareness

AI Engine

Control Logic

Actuator Output Interface

Actuators (examples)

Remote management

Firewall and routing rules

Device reset

Chat message

Report

# Case study: A user mistakenly downloads and executes a malicious file

# Case study

| Initial hardening | Monitoring & response | Reporting | Recovery |

- Configure all clients to send HTTP(s) traffic via a proxy

- Enable detailed logging and send logs to a central server

- Set up recording of all network traffic and feature extraction

# Case study

| Initial hardening | Monitoring & response | Reporting | Recovery |

- Proxy detects the malicious payload

- Logging reports the execution of an unknown file

- Sniffer detects connection to C&C server

# Case study

| Initial hardening | Monitoring & response | Reporting | Recovery |
|---|---|---|---|

- Proxy detects the malicious payload → Remove payload

- Logging reports the execution of an unknown file → Block execution

- Sniffer detects connection to C&C server → Drop packets

# Case study

| Initial hardening | Monitoring & response | **Reporting** | Recovery |

- Human-readable report with information about the incident
  - Malware source (compromised webserver)
  - Malware type
  - …

# Case study

| Initial hardening | Monitoring & response | Reporting | **Recovery** |

- Restore a device if the malware was executed

  E.g., from a backup

# Current status of the project

Sensors (examples)

Actuators (examples)

```
Network traffic  →  Sensor Input Interface  →  Situational Awareness  ↔  AI Engine
Event logs                                                            ↔  Control Logic  →  Actuator Output Interface
External monitoring
Support tickets
Device credentials
```

Network traffic

Event logs

External monitoring

Support tickets

Device credentials

Sensor Input Interface

Situational Awareness

AI Engine

Control Logic

Actuator Output Interface

Remote management

Firewall and routing rules

Device reset

Chat message

Report

# Current status of the project

Sensors (examples)

Actuators (examples)

Network traffic

Event logs

External monitoring

Support tickets

Device credentials

Sensor Input Interface

Situational Awareness

AI Engine

Control Logic

Actuator Output Interface

Remote management

Firewall and routing rules

Device reset

Chat message

Report

# Training and testing AI models requires high-quality datasets

- We work with data from the Swiss Blue Team and other collaborating nations

- Since 2023, we (researchers) participate as a separate Blue Team in Locked Shields for data collection

- Datasets collected in 2023 and 2024 are (soon) publicly available

# The LSPR23 dataset

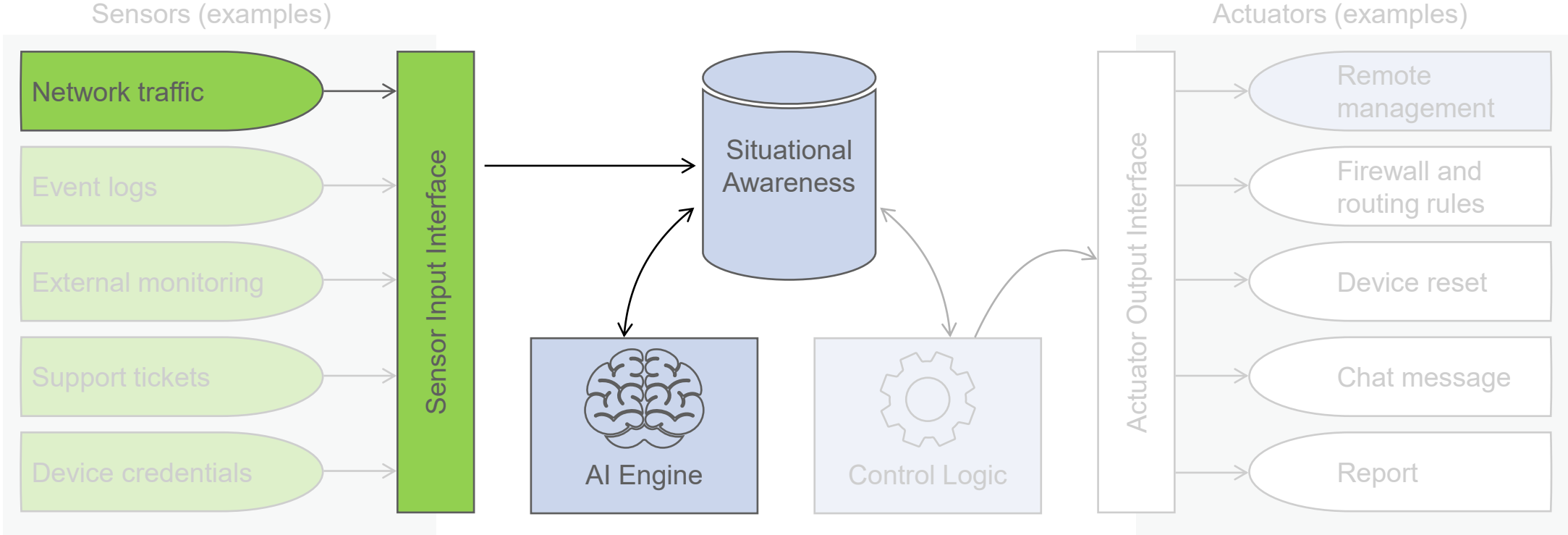| Total Flows | | 16,353,511 |
|---|---|---|
| Labels: | Source | Destination |
| Blue Flows | 12,493,826 | 10,821,533 |
| Red Flows | 13,880 | 1,630,719 |
| Green Flows (not scoring) | 17,256 | 263,318 |
| Green Flows (scoring bots) | 3,415,280 | 136,321 |
| External Flows | 413,269 | 3,501,620 |
| Benign Flows* | | 14,708,912 |
| Malicious Flows* | | 1,644,599 |
| Network segments: | | |
| Berilia Energy Group | 7,566,863 | 6,755,988 |
| Berilia Airforce | 4,573,708 | 3,849,478 |
| Bank Of Berilia | 277,705 | 188,984 |
| Berilia Airforce 5G | 68,369 | 27,649 |
| Other | 3,866,866 | 5,531,412 |

| Attack labels "Goals" | |
|---|---|
| Privilege Escalation | 63 |
| System Compromise | 58 |
| Data Theft | 52 |
| Website Defacement | 47 |
| Non Destructive | 18 |
| Attack labels "Methods" | |
| Remote Code Execution | 35 |
| Authorize with Default credentials | 28 |
| Remote Desktop | 23 |
| Authorize with RT credentials | 10 |
| Authorize with Stolen credentials | 9 |

# Current status of the project

Network traffic

Event logs

External monitoring

Support tickets

Device credentials

Sensor Input Interface

Situational Awareness

AI Engine

Control Logic

Actuator Output Interface

Remote management

Firewall and routing rules

Device reset

Chat message

Report

# We used supervised learning to detect "Command and Control" traffic

- Initial results showed that the approach works well in some cases (when training and testing was in similar settings), but it does not generalize well

Test data

| Training data | | LS17 | LS18 |
|---|---|---|---|
| | LS17 | 0.993 | 0.966 |
| | LS18 | 0.945 | 0.993 |
| | LS19 | 0.743 | 0.928 |
| | LS21 | 0.952 | 0.918 |

F1 scores

# We used supervised learning to detect "Command and Control" traffic

- Initial results showed that the approach works well in some cases (when training and testing was in similar settings), but it does not generalize well

Test data

| Training data | | LS17 A | LS18 A | LS19 A | LS21 A | LS21 B |
|---|---|---|---|---|---|---|
| | LS17 A | 0.993 | 0.966 | 0.007 | 0.856 | 0.215 |
| | LS18 A | 0.945 | 0.993 | 0.060 | 0.806 | 0.167 |
| | LS19 A | 0.743 | 0.928 | 0.791 | 0.351 | 0.000 |
| | LS21 A | 0.952 | 0.918 | 0.038 | 0.986 | 0.158 |

F1 scores

# Current status of the project

# We detect and fix common misconfigurations automatically

- Currently two tools:

  - Automatically change all login credentials

  - Scan the network for nginx web servers, analyze their

    configuration and fix some misconfigurations
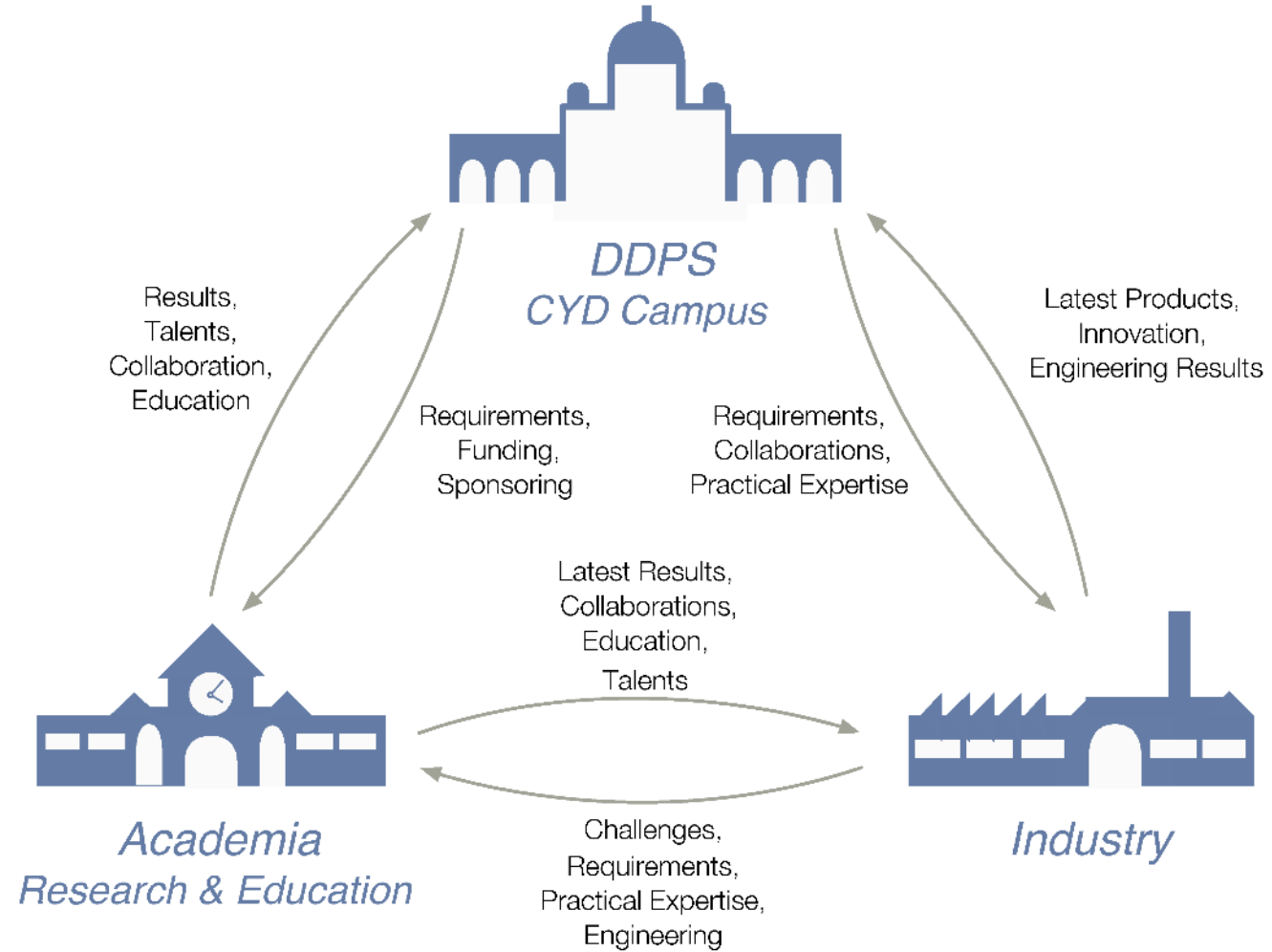
    automatically

# Other topics we are currently (planning to) investigate

- Using Generative AI / LLMs

  - Parse support tickets

  - Generate reports

  - Generate code or configuration

- Improve existing models such that they generalize better

- Building an additional training and testing environment

# The Cyber-Defence Campus connects government, academia and industry

**Thank you for your attention!**

Dr. Roland Meier

roland.meier@ar.admin.ch

cydcampus.admin.ch

cyber-defence-campus