Towards an Al-powered Player in Cyber Defence Exercises

Roland Meier⁽¹⁾, Artūrs Lavrenovs⁽²⁾, Kimmo Heinäaro⁽²⁾,

Luca Gambazzi⁽³⁾, Vincent Lenders⁽³⁾

26.05.2021









(3)

2'500 Attacks



2'500 Attacks

against

4'000 Systems



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

Picture: NATO CCDCOE

11:30:42

CCDCOE

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

Red Team vs. Blue Team exercise Attackers Defenders 1 Team 1 Team / country





2 years ago, we presented a system which uses AI to identify C&C channels

2019 11th International Conference on Cyber Conflict: Silent Battle T. Minárik, S. Alatalu, S. Biondi, M. Signoretti, I. Tolga, G. Visky (Eds.) 2019 © NATO CCD COE Publications, Tallinn Permission to make digital or hard copies of this publication for internal use within NATO and for personal or educational use when for non-profit or non-commercial purposes is granted providing that copies bear this notice and a full citation on the first page. Any other reproduction or transmission requires prior written permission by NATO CCD COE.

Machine Learning-based Detection of C&C Channels with a Focus on the Locked Shields Cyber Defense Exercise

Nicolas Känzig

Department of Information Technology and Electrical Engineering ETH Zürich Zürich, Switzerland kaenzign@student.ethz.ch

Luca Gambazzi

Roland Meier

Department of Information Technology and Electrical Engineering ETH Zürich Zürich, Switzerland meierrol@ethz.ch

Vincent Lenders



What would it take to have a fully automated Blue Team in a future iteration of Locked Shields?



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





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



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 ticket	
Yellow Team	Provide periodic reports	
White team	Voice or video calls through the Gamenet	
Green Team	Gamenet status, reverting of devices	





Initial hardening

Prepare systems before attacks start







Initial hardening	Monitoring a	& response	Repor	ting
Prepare systems before attacks start			Respond to u and provide in	iser requests icident reports
	Monitor sys they are ur and mitiga	stems while nder attack ate attacks		



Initial hardening	Monitoring & response	Reporting	Recove	ery
Prepare systems before attacks start		Respond to user request and provide incident repo	s rts	
	Monitor systems while they are under attack and mitigate attacks		Restore a compromise	nd fix d systems



























Organizational inputs



















Organisational inputs define the environment and the objectives

System documentation

Devices and services running in the Gamenet



Credentials

To access and configure devices

Scoring system

Defines the objectives of the Blue Team

External monitoring

Information about the current performance



Application-level sensors provide information about running applications



- Application configuration
 E.g., users and permissions
- Application state
 - E.g., update status, access logs



Device-level sensors report the status of all clients and servers



System information

E.g., running services, update status

System logs

E.g., login attempts, executed commands

Active monitoring

E.g., for additional logs, screen capture



Network-level sensors provide access to network traffic



- Network configuration
 E.g., topology, open firewall ports
- Network state
 E.g., current load, reachability
- Network traffic
 Capture for inspection



User-interaction sensors receive chat and email messages

Support tickets

Notifications about observed problems





Overview





There are four points where actuators can act





Application-level actuators



Application firewall

E.g., block malicious uploads in web application firewall

- Application configuration
 E.g., restrict access
- Application patching
 E.g., install updates



Device-level actuators



Remote management

Log in to every device and run commands

Restore initial state

Revert a device to a previous state (e.g., from a backup)



Network-level actuators

Network configuration
 E.g., firewall rules





User-interaction actuators



Chat and e-mail

Respond to messages in a human-readable format

Incident reports

Report successful and failed attacks to the Green Team



Overview





The situational awareness database contains all the relevant information

- Sensor data
 - E.g., access logs
- Al models
 - E.g., anomaly detection
- Actuator outputs
 E.g., incident reports
- External information
 - E.g., known software vulnerabilities



Overview









Level 0: Reactive narrow AI





Level 1: Limited-memory narrow AI

"Machine learning" today

Level 0: Reactive narrow AI





Level 2: General AI

Mimics human intelligence

Level 1: Limited-memory narrow AI

Level 0: Reactive narrow AI

"Machine learning" today





Level 3: Super Al

Level 2: General AI

Surpasses human intelligence

Mimics human intelligence

Level 1: Limited-memory narrow AI

"Machine learning" today

Level 0: Reactive narrow AI



Only level 0 and level 1 exist today



Level 3: Super Al

Level 2: General Al

Surpasses human intelligence

Mimics human intelligence

Level 1: Limited-memory narrow AI

Level 0: Reactive narrow AI

"Machine learning" today



Five tasks for AI

Identification / classification

Categorisation

What is it?

What belongs together?

Assessment

Recommendation

What is important?

What should be done?

Prediction

What will happen?



AI for initial hardening

Initial hardening	Monitoring & response	Reporting		Recovery
Identifi	cation / classification			
Categor	risation	Find groups of similar	devices	
Assessr	nent			
Recomm	nendation			
Predicti	ion			



AI for monitoring and response

Initial hardening	Monitoring & response	Reporting	Recovery
Ident	ification / classification	Detect malicious netwo	ork traffic
Cate	gorisation	Detect malicious patte	rns in log files
Asses	ssment		
Reco	mmendation		
Predi	iction		





Initial hardening	Monitoring & response	Reporting		Recovery
Identif	ication / classification			
Catego	prisation	Link support tickets and	d mon	itoring alerts
Assess	ment	Prioritise support ticke	ts	
Recom	mendation	Formulate response to	suppo	ort tickets
Predic	tion	Predict impact on the s	scoring	5





Initial hardening	Monitoring & response	Reporting	Recovery
Identifi	cation / classification	Find devices that need	to be recovered
Catego	risation	Find similar devices as	a template for recovery
Assessr	nent		
Recom	mendation		
Predict	ion		



Overview







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



Actions for initial hardening

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

Software deployment

E.g., web application firewall, antivirus, monitoring agents

Patch systems

E.g., install updates, find backdoors

Secure initial configuration

E.g., change passwords, add firewall rules



Actions for monitoring and response

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

- Disable exploited services
 - E.g., compromised websites
- Monitor and block malicious system calls
 E.g., process creations or opening of files
- Monitor and block malicious network connections
 E.g., command and control flows



Actions for reporting

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

Generate human-readable reports

E.g., about mitigated and successful attacks



Actions for recovery

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

- Self-recovery from previous backups
 - E.g., database snapshots
- Revert devices to their initial state

E.g., special systems without backup possibilities



Overview





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







Initial hardening	Monitoring & response	Reporting	Recovery	

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

「「」「二
<u>ک</u>
T

- Enable detailed logging and send logs to a central server
- Set up recording of all network traffic and feature extraction





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





Initial hardening Monitoring & response Reporting Recovery
--



Proxy detects the malicious payload

 \rightarrow Remove payload

- Logging reports the execution of an unknown file
- Sniffer detects connection to C&C server

- \rightarrow Block execution
- \rightarrow Drop packets





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

C	
	النيريز
	്പ∥
C	T

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

...





Restore a device if the malware was executed



E.g., from a backup



What comes next?



- Implementation
- Evaluation in Locked Shields 20xx
- Adaptation for other exercises



What comes next?



- Implementation
- Evaluation in Locked Shields 20xx
- Adaptation for other exercises

Questions?



Roland Meier meierrol@ethz.ch