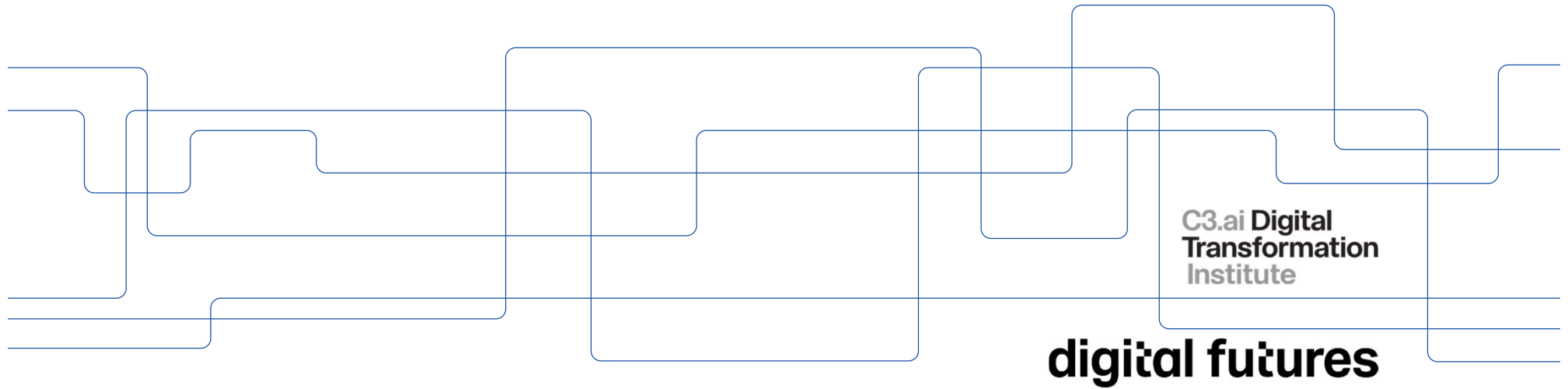# Boosting Cyber Resilience with Human-in-the-loop AI
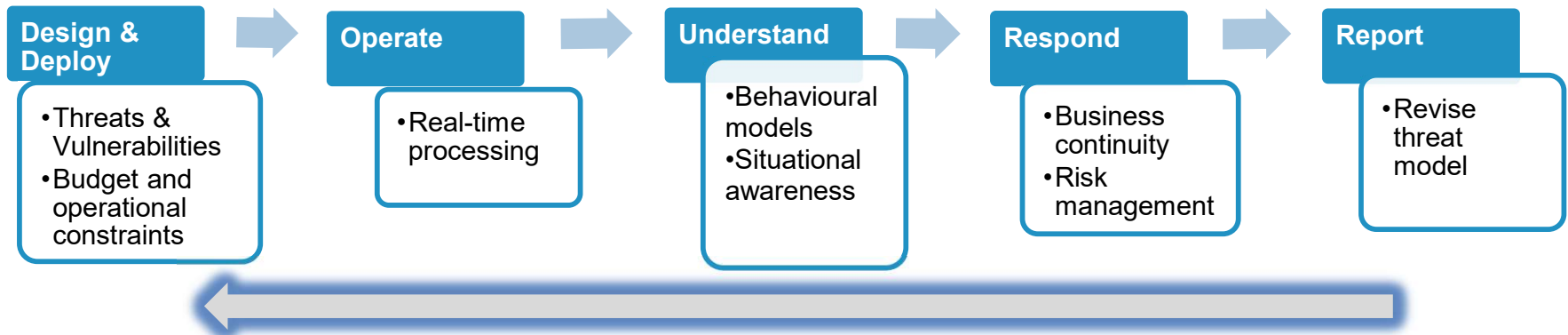
György Dán

IEEE CNS 2024 Workshop on Cyber Resilience

# The Needle in the Haystack?





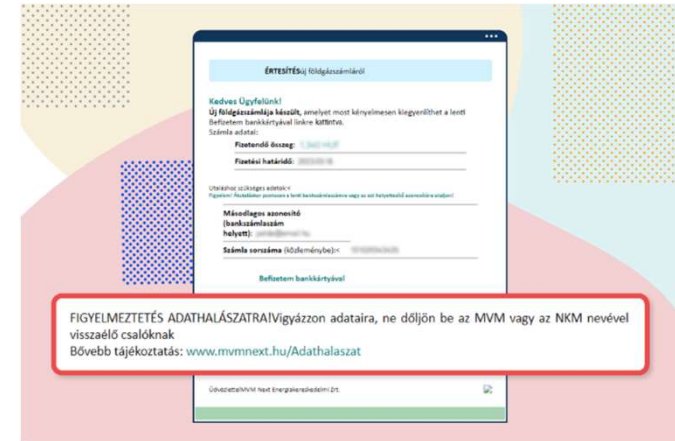| Design & Deploy | | Operate | | Understand | | Respond | | Report |
|---|---|---|---|---|---|---|---|---|
| •Threats & Vulnerabilities<br>•Budget and operational constraints | | •Real-time processing | | •Behavioural models<br>•Situational awareness | | •Business continuity<br>•Risk management | | •Revise threat model |

# AI Changes the Threat Model

# AI-Powered Adversaries

- Social engineering
  - Target selection, deepfakes

- Phishing
  - Improved personalization, live communication at scale

- Vulnerability discovery
  - Hardware/software vulnerability analysis

- Autonomous malware



## Lore a Red Team Emulation Tool

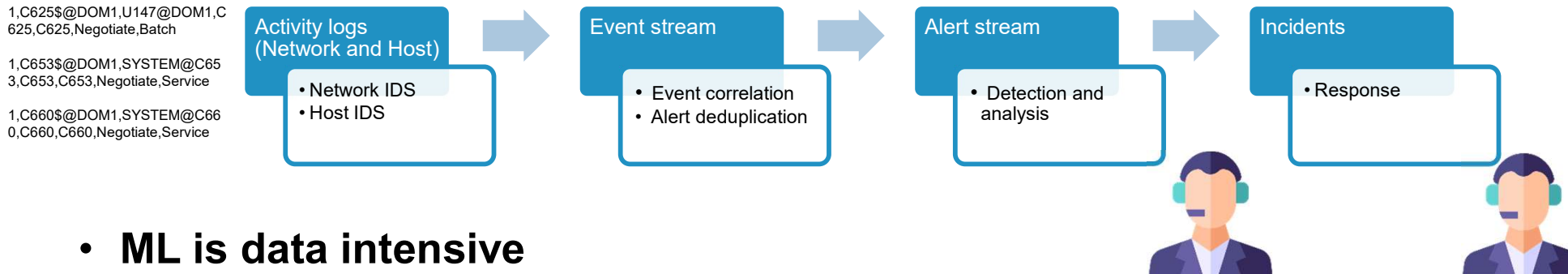Publisher: IEEE | Cite This | PDF



## New type of polymorphic fully autonomous malware uses AI
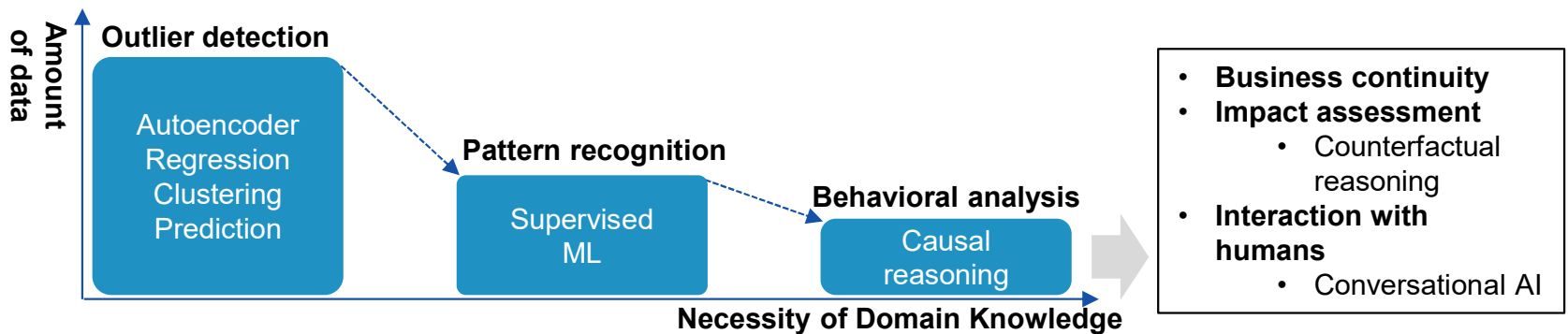
Technology News | August 2, 2023

# AI-Powered Cyber Resilience

- **From logs to incident response**

```
1,C625$@DOM1,U147@DOM1,C
625,C625,Negotiate,Batch

1,C653$@DOM1,SYSTEM@C65
3,C653,C653,Negotiate,Service

1,C660$@DOM1,SYSTEM@C66
0,C660,C660,Negotiate,Service
```

**Activity logs (Network and Host)**
- Network IDS
- Host IDS

**Event stream**
- Event correlation
- Alert deduplication

**Alert stream**
- Detection and analysis

**Incidents**
- Response

- **ML is data intensive**

Amount of data

**Outlier detection**

Autoencoder
Regression
Clustering
Prediction

**Pattern recognition**

Supervised ML

**Behavioral analysis**

Causal reasoning

**Necessity of Domain Knowledge**

- **Business continuity**
- **Impact assessment**
  - Counterfactual reasoning
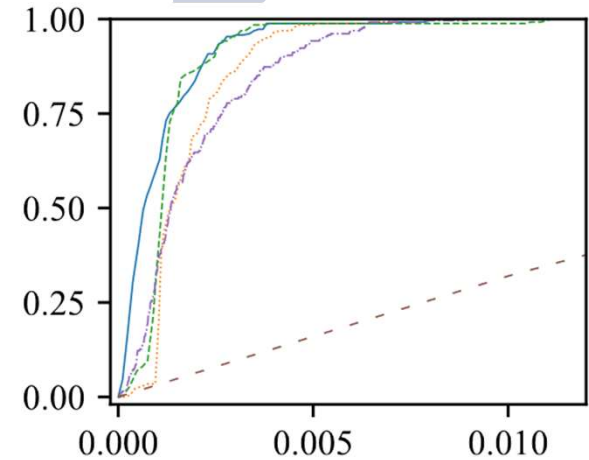- **Interaction with humans**
  - Conversational AI

# ML/AI as a Power Tool



```
Feb 10 15:45:09 ubuntu-lts sshd[47341]: Failed password for root from
103.106.189.143 port 60824 ssh2
Feb 10 15:45:11 ubuntu-lts sshd[47341]: Connection closed by authenticating user
root 103.106.189.143 port 60824 [preauth]
Feb 10 15:45:11 ubuntu-lts sshd[47339]: Failed password for root from
180.101.88.228 port 11349 ssh2
Feb 10 15:45:12 ubuntu-lts sshd[47343]: pam_unix(sshd:auth): authentication
failure; logname= uid=0 euid=0 tty=ssh ruser= rhost=103.106.189.143   user=root
Feb 10 15:45:14 ubuntu-lts sshd[47339]: Failed password for root from
180.101.88.228 port 11349 ssh2
Feb 10 15:45:14 ubuntu-lts sshd[47343]: Failed password for root from
103.106.189.143 port 33990 ssh2
Feb 10 15:45:16 ubuntu-lts sshd[47343]: Connection closed by authenticating user
root 103.106.189.143 port 33990 [preauth]
```

Tokenizer → Bidirectional Transformer → Prediction loss

Gökstorp et al,``Anomaly Detection in Security Logs using Sequence Modeling,'' in Proc. of IFIP/IEEE NOMS, 2024
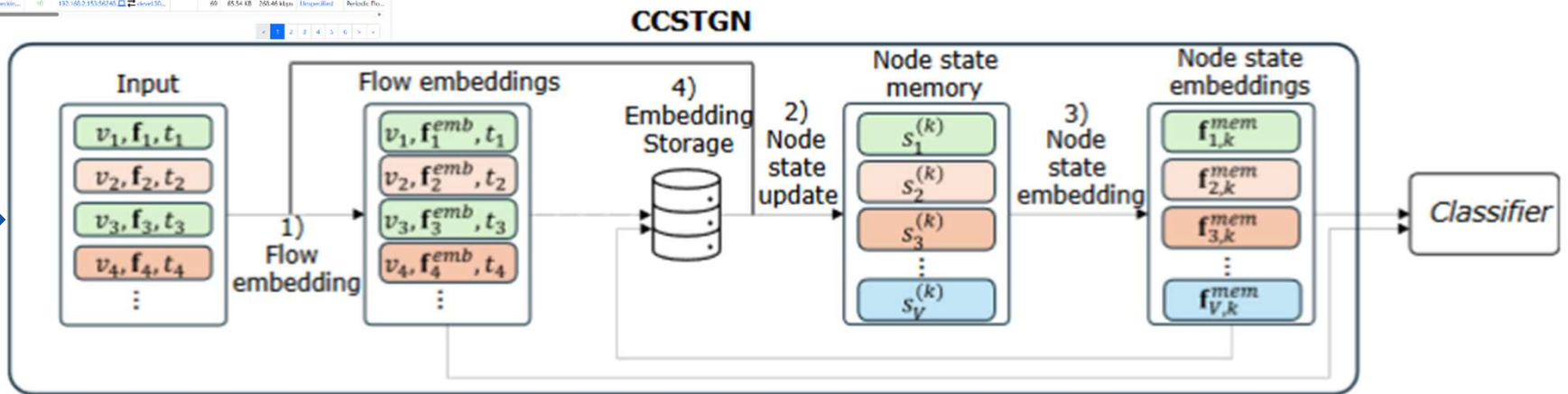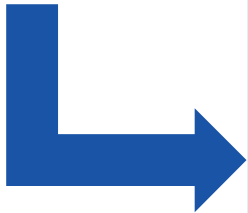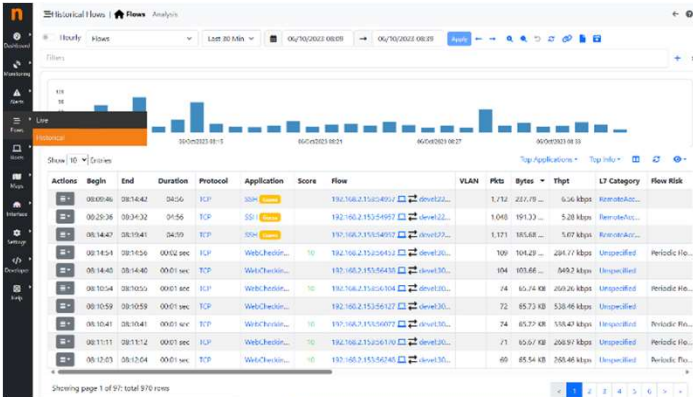
# ML/AI as a Power Tool



Rehman et al., "FLASH: A Comprehensive Approach to Intrusion Detection via Provenance Graph Representation Learning", in Proc. of IEEE S&P, 2024
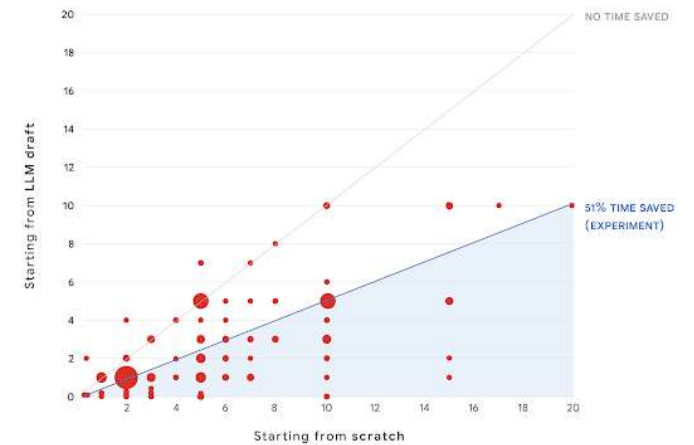
# ML/AI as a Power Tool



Santos et al., ``Channel-Centric Spatio-Temporal Graph Networks for Network-based Intrusion Detection,'' in Proc. of IEEE CNS, 2024

# ML/AI as a Power Tool
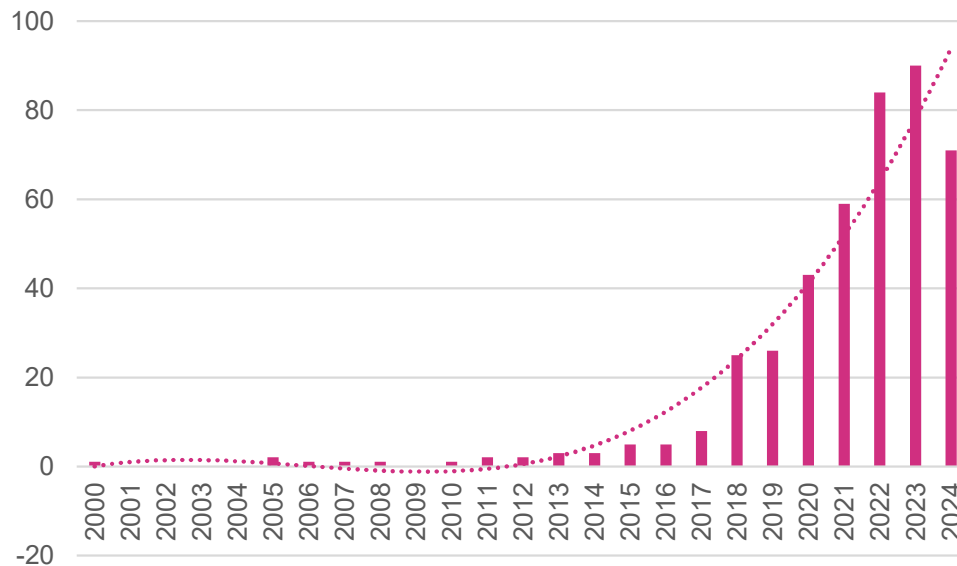


```
<Security Incident>
<Title> [tool_name_verdict] Abuse verdict for project id: xyz.</Title>
<Metadata> This ticket was filled and submitted on the 2023-10-01. It was marked with the labels:
"Investigation" and "AB".</Metadata>
<Description> Counter-Abuse has issued an abuse verdict against a GCP project.</Description>
<Additional Information> The incident was reported through the xyz pipeline with a policy violation
of "COIN_MINING".
The infraction can be found in the project xyz.</Additional Information>
<Date Incident> 2023-10-01 11:50:19</Date Incident><Incident Causes> The identified causes are:
MISCONFIGURATION, WEAK_OR_NO_PASSWORD</Incident Causes><Actions Taken> The following actions were
taken:
1) Action1
2) Action2</Actions Taken>
<Software Involved> Software1</Software Involved>
<Sensitive Data> - NONE, TEST</Sensitive Data>
<Mitigation History><Comment index="1" author="user1@domain.com"> Looks like there was a CPU spike:
URL around 05:00. Running application1 now.</Comment>
<Comment index="2" author="user3@domain.com"> Instance compromised, shutting it down</Comment>
<Comment index="3" author="user4@domain.com"> InstanceMetadata</Comment>
<Comment index="4" author="user@domain.com"> Get additional information on InstanceMetadata:
URL`<Code Section/>`</Comment>
<Comment index="5" author="user3@domain.com"> Looks like it was compromised through  successfully
authentication as root account using SSH with password authentication: `<Code Section/>`</Comment>
<Comment index="6" author="user3@domain.com"> A malicious cron job was created on the machine
`<Code Section/>`. The cron job downloaded a bash script from IP and executed it. The script was
not present under `<Code Section/>` at the time of the investigation `<Code Section/>`</Comment>
<Comment index="10" author="user3@domain.com"> Exec update sent.</Comment>
</Mitigation History>
```

Time spent (in minutes) writing an incident summary

NO TIME SAVED

51% TIME SAVED
(EXPERIMENT)

Starting from LLM draft

Starting from scratch

– https://security.googleblog.com/2024/04/accelerating-incident-response-using.html

# Autonomous Cyber Defense



Introduction of AI/ML

- **Operational requirements**
  - Auditable
  - Controllable
  - Transferable/Adaptive
  - Secure
  - Observable/Explainble

Source: https://github.com/Limmen/awesome-rl-for-cybersecurity

# Human-in-the-loop AI for Security

- **AI/ML complementing human decision making**
  - Reduced response time
  - Higher accuracy



**Intelligent Vehicle**

**Robot**

Human Model — Robot Model — Task Model

Driver

Perception / Action

Human (user/collaborator)

Perception / Action

Task/Environment

Environment

https://thebossmagazine.com

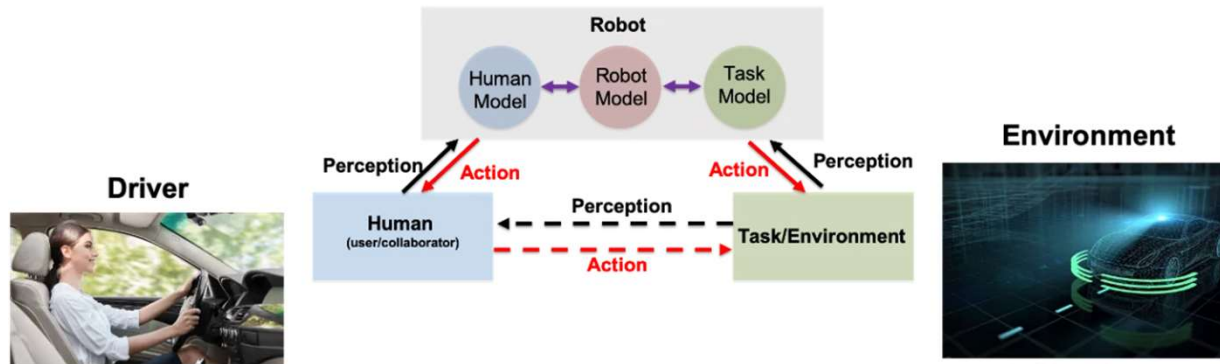https://arm.stanford.edu/research/leveraging-human-intent-shared-autonomy
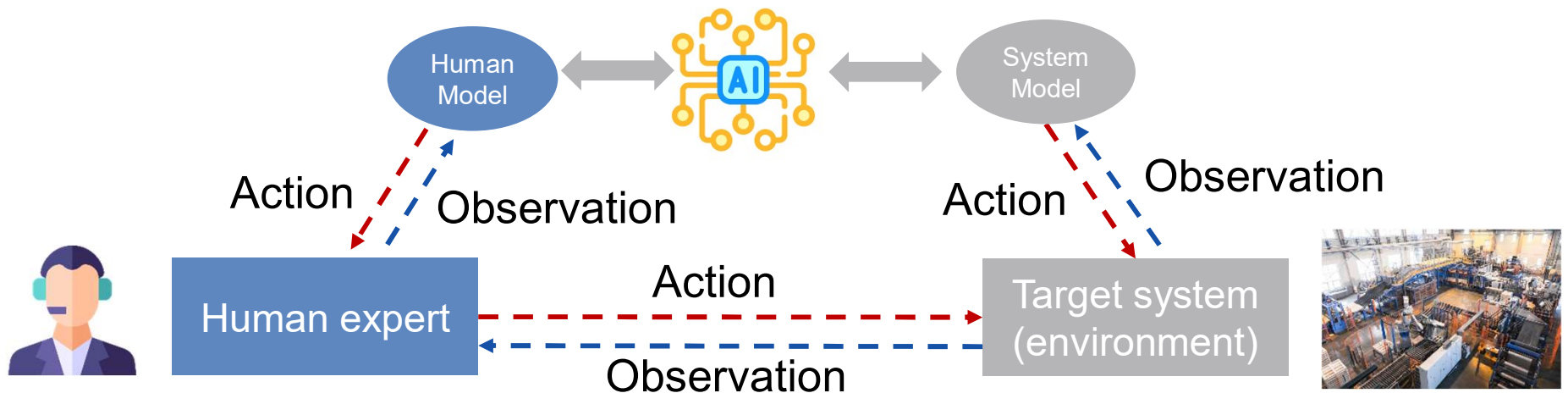
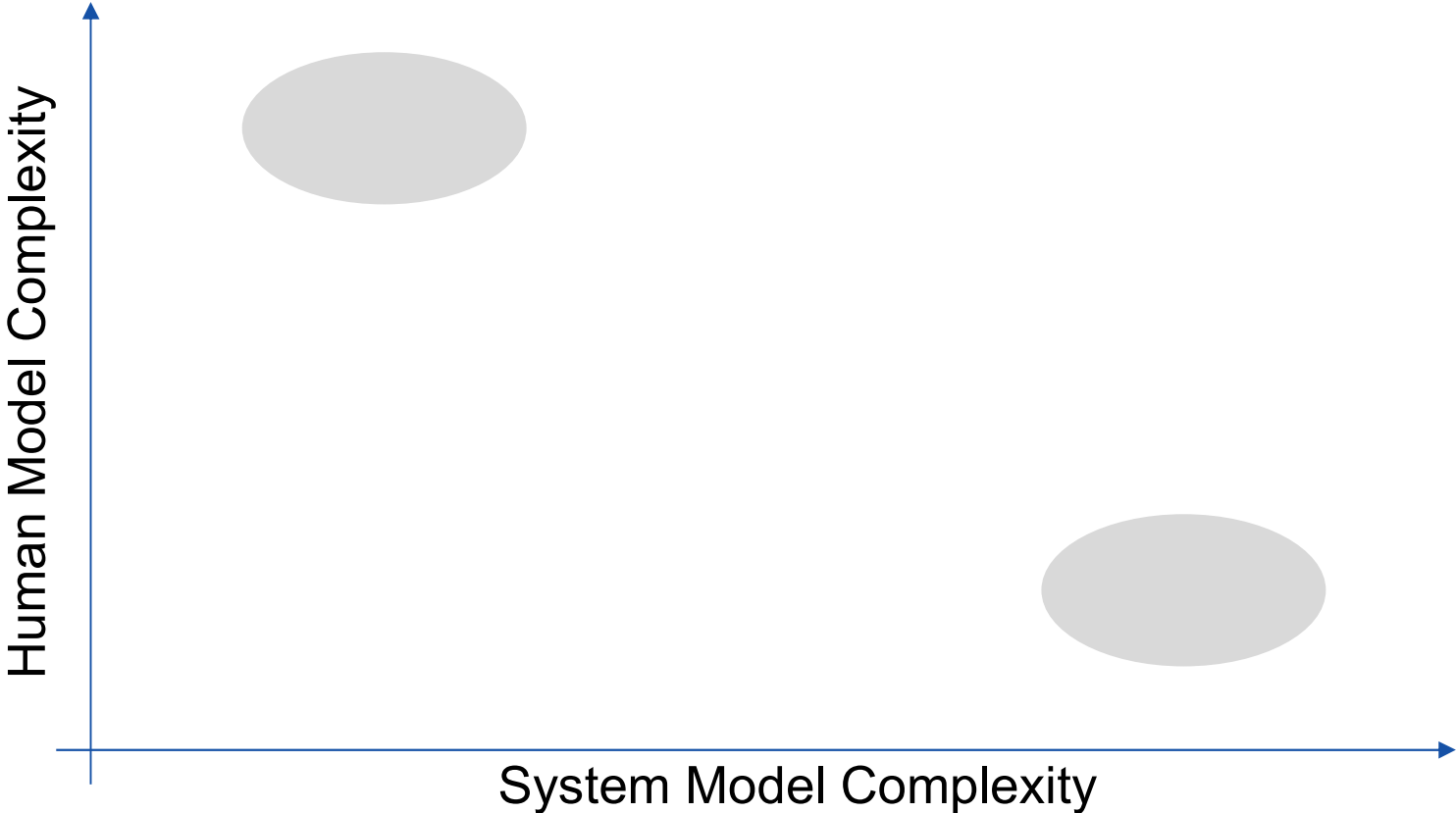# Human-in-the-loop AI for Security

- **AI/ML complementing human decision making**
  - Reduced response time
  - Higher accuracy

# Framework Design Space

# Human-in-the-Loop AI Framework

# State in Cyber Security

- Attack tree: Hypergraph of conditions and exploits

- Attack state: the set of conditions/privileges the attacker gained



- States and transitions → Markov model

# Problem of Partial Observability

- Security state is not visible to the defender

  - Attacker activity can trigger alerts



$X$ — states
$y$ — possible observations
$a$ — state transition probabilities
$b$ — output probabilities

- Hidden Markov model

# System model – Security state

- **Time is slotted**

- **Attack Hypergraph**
  - Nodes: conditions (access privilege, etc)
  - Hyperedges: exploits

$\mathcal{H} = (\mathcal{N}, \mathcal{E})$ where $\mathcal{N} = \{c_1, ..., c_{n_c}\}$, and $\mathcal{E} = \{e_1, ..., e_{n_e}\}$
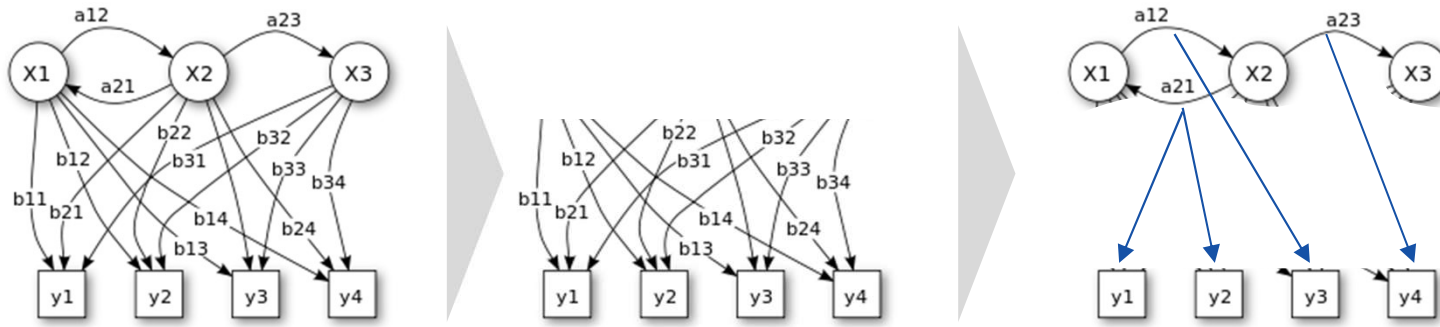
- **Security state: set of enabled conditions**

$$s_1 = \{c_1\} \qquad s_2 = \{c_1, c_2\}$$

- **Example**
  - c1: wu-ftpd 2.5 running on host
  - c2: ftp server remotely accessible
  - e3: CVE-1999-0878
  - c3: Root privilege on host



Kim et al, "An Active Learning Approach to Dynamic Alert Prioritization for Real-time Situational Awareness" Proc. of IEEE CNS, 2022

# Attacker model

- **Attacker chooses exploits independently**
  - Probability of choosing exploit $e_i$: $\alpha_{e_i}$
  - Probability that exploit $e_i$ succeeds: $\beta_{e_i}$

- **If exploit $e_i$ is used**
  - Generates alert $a$ with probability $\delta_{ia}$

- **False positive** with probability $\zeta_a$

- **Alert vector** $Y_t = \left(y_1, \ldots, y_{n_z}\right)$



Kim et al, "An Active Learning Approach to Dynamic Alert Prioritization for Real-time Situational Awareness" Proc. of IEEE CNS, 2022

# Defender model

- **Observation at time $t$:** $Y_t$ (alert vector)

- **Action:**
  - Inspect up to $I$ alerts in $Y_t$
  - Inspecting alert $y_t^a$ results in modified alert $\hat{y}_t^a$

- **Human model:** Investigation error probability $\omega$
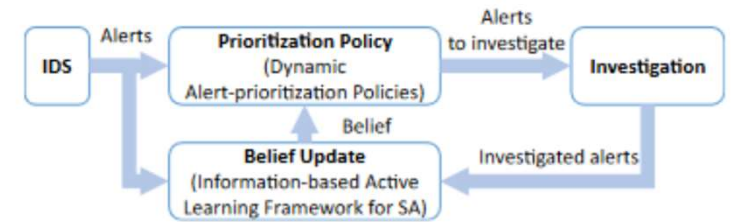
- **Belief about security state**

$$\pi_t = \begin{bmatrix} \pi_t^{1,1} & \pi_t^{1,2} & \cdots & \pi_t^{1,n_a} \\ \pi_t^{2,1} & \pi_t^{2,2} & \cdots & \pi_t^{2,n_a} \\ \vdots & \vdots & \ddots & \vdots \\ \pi_t^{n_s,1} & \pi_t^{n_s,2} & \cdots & \pi_t^{n_s,n_a} \end{bmatrix} \in \Delta(\mathcal{S} \times \Phi)$$



| | Investigation outcome | |
|---|---|---|
| Ground truth | TP | FP |
| TP | $1 - \omega$ | $\omega$ |
| FP | $\omega$ | $1 - \omega$ |

- **Cost:** State estimation error $\quad J^\kappa = \lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \gamma^t MSE(\pi_t^\kappa, s_t^\kappa)$

- **Optimal policy:** $\kappa^* \in \arg\min_{\kappa \in \mathcal{K}} J^\kappa$

# Active learning for alert prioritization



- **In practice the state is unknown → cannot calculate MSE**

  > *Use **belief uncertainty as a proxy** for the MSE.*
  > *Intuition: Low uncertainty is likely to imply an accurate belief*

- **Proposed candidate policies**

  – Max-entropy

  → Investigate the alert $v$ that decreases the entropy most
  $$\min_{v} H(S_{t+1} = s_{i'}, \Phi_{t+1} = \phi_{l'} | V_{t+1} = v, Y_{t+1} = y_n, \Pi_t = \pi_t)$$

  – Bayes factor policy

  → Investigate the most ambiguous alert
  (alert probability without false positives vs. false positive rate)
  $$K^a = \frac{P(Y_{t+1}^a = 1 \mid Y_{t+1}^{-a} = y_n^{-a}, \Pi_t = \pi_t)|_{\zeta_a=0}}{\zeta_a},$$



Kim et al, "An Active Learning Approach to Dynamic Alert Prioritization for Real-time Situational Awareness" Proc. of *IEEE CNS,* 2022

# System Level Benefit



Less human effort needed

Less skilled experts needed

Kim et al, "An Active Learning Approach to Dynamic Alert Prioritization for Real-time Situational Awareness" Proc. of *IEEE CNS,* 2022

# Framework Design Space



Compressed CoreLang Attack Graph with 540 nodes



Human Model Complexity

System Model Complexity

Katsikeas et al. "An attack simulation language for the IT domain," in *Proc. of Int. Workshop on Graphical Models for Security*, pp. 67–86, 2020.

# MITRE ATT&CK Model



| Initial Access | Execution | Persistence | Evasion | Discovery | Lateral Movement | Collection | Command and Control | Inhibit Response Function | Impair Process Control | Impact |
|---|---|---|---|---|---|---|---|---|---|---|
| Data Historian Compromise | Change Program State | Hooking | Exploitation for Evasion | Control Device Identification | Default Credentials | Automated Collection | Commonly Used Port | Activate Firmware Update Mode | Brute Force I/O | Damage to Property |
| Drive-by Compromise | Command-Line Interface | Module Firmware | Indicator Removal on Host | I/O Module Discovery | Exploitation of Remote Services | Data from Information Repositories | Connection Proxy | Alarm Suppression | Change Program State | Denial of Control |
| Engineering Workstation Compromise | Execution through API | Program Download | Masquerading | Network Connection Enumeration | External Remote Services | Detect Operating Mode | Standard Application Layer Protocol | Block Command Message | Masquerading | Denial of View |
| Exploit Public-Facing Application | Graphical User Interface | Project File Injection | Rogue Master Device | Network Service Scanning | Program Organization Units | Detect Program State | | Block Reporting Message | Modify Control Logic | Loss of Availability |
| External Remote Services | Man in the Middle | System Firmware | Rootkit | Network Sniffing | Remote File Copy | I/O Image | | Block Serial COM | Modify Parameter | Loss of Control |
| Internet Accessible Device | Program Organization Units | Valid Accounts | Spoof Reporting Message | Remote System Discovery | Valid Accounts | Location Identification | | Data Destruction | Module Firmware | Loss of Productivity and Revenue |
| Replication Through Removable Media | Project File Injection | | Utilize/Change Operating Mode | Serial Connection Enumeration | | Monitor Process State | | Denial of Service | Program Download | Loss of Safety |
| Spearphishing Attachment | Scripting | | | | | Point & Tag Identification | | Device Restart/Shutdown | Rogue Master Device | Loss of View |
| Supply Chain Compromise | User Execution | | | | | Program Upload | | Manipulate I/O Image | Service Stop | Manipulation of Control |
| Wireless Compromise | | | | | | Role Identification | | Modify Alarm Settings | Spoof Reporting Message | Manipulation of View |
| | | | | | | Screen Capture | | Modify Control Logic | Unauthorized Command Message | Theft of Operational Information |
| | | | | | | | | Program Download | | |
| | | | | | | | | Rootkit | | ATT&CK for Enterprise |
| | | | | | | | | System Firmware | | ATT&CK for ICSs |
| | | | | | | | | Utilize/Change Operating Mode | | |

Choi et al., "Probabilistic Attack Sequence Generation and Execution Based on MITRE ATT&CK for ICS Datasets", in Proc. of ACM CSET, 2021
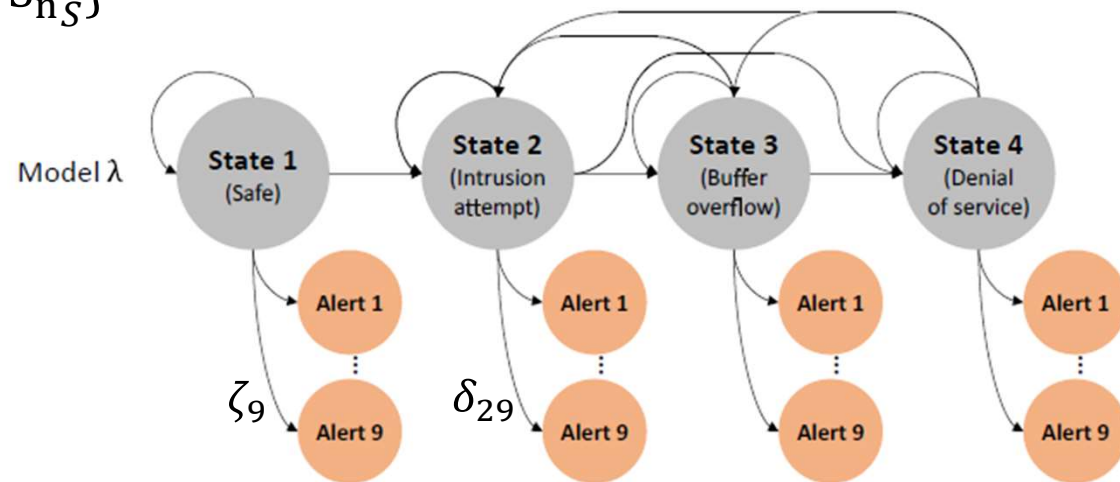
# Attack and Observation Model

- Set of attacker states $\mathcal{S} = \{s_1, \dots, S_{n_S}\}$
- State at time $t$: $S_t$
- Set of alerts $\mathcal{J} = \{1, \dots, J\}$

- True alert probability
$$\delta_{ij} = P(Y_t^j = 1 | S_t = s_i)$$

- False alert probability
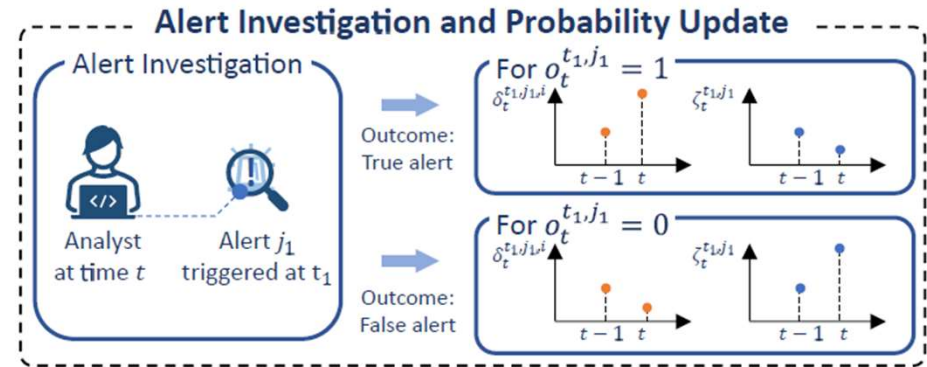$$\zeta_j = P(Y_t^j = 1 | S_t = s_1)$$
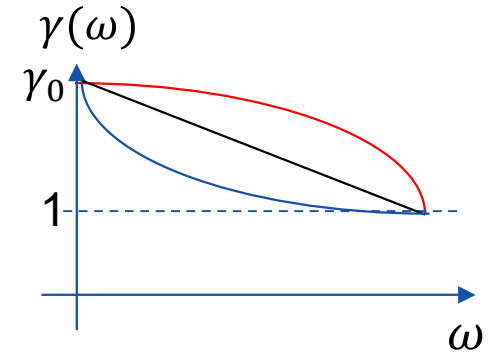
# Defender Model

- **Observes** alerts $Y_t$ at time $t$

- **Investigates** up to $I$ alerts $v \subseteq Y_{1:t}$
  - Investigation outcome $o_t$

- **Human model:** Investigation error probability $\omega$

- Confidence function

$$\gamma(\omega) = \begin{cases} 2(1-\gamma_0)\omega + \gamma_0, & (linear), \\ 4(1-\gamma_0)\omega^2 + \gamma_0, & (concave), \\ 4(\gamma_0-1)(\omega-0.5)^2 + 1 & (convex), \end{cases}$$

- Update of HMM Observation Model

$$\delta_t^{t',j,i} = \begin{cases} \frac{1}{\gamma(\omega)} \delta_{t-1}^{t',j,i} & if \ o_t^{t',j} = 0 \\ min\left(\gamma(\omega)\delta_{t-1}^{t',j,i}, 1\right) & if \ o_t^{t',j} = 1 \end{cases}$$

$$\zeta_t^{t',j} = \begin{cases} min\left(\gamma(\omega)\zeta_{t-1}^{t',j}, 1\right) & if \ o_t^{t',j} = 0 \\ \frac{1}{\gamma(\omega)} \zeta_{t-1}^{t',j} & if \ o_t^{t',j} = 1 \end{cases}$$





Alert Investigation and Probability Update

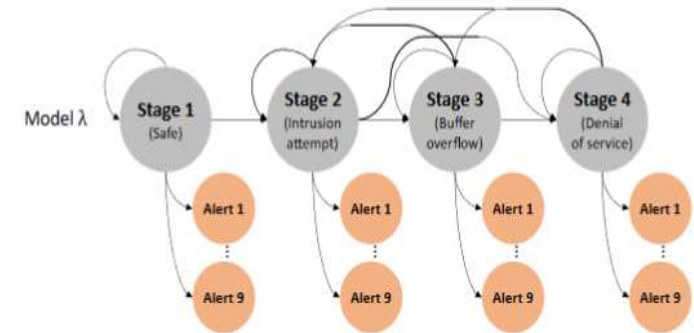Kim et al, "Human-in-the-loop Cyber Intrusion Detection Using Active Learning" IEEE TIFS, 2024

# Defender's Problem

- **Defender objective:** Minimize mean time to detection

$$\kappa^* = \arg\min_{\kappa \in \mathcal{K}} \sup_{t_{1 \to 2} > 0} \mathbb{E}^{(t_{1 \to 2})}[d^\kappa - t_{1 \to 2}]$$

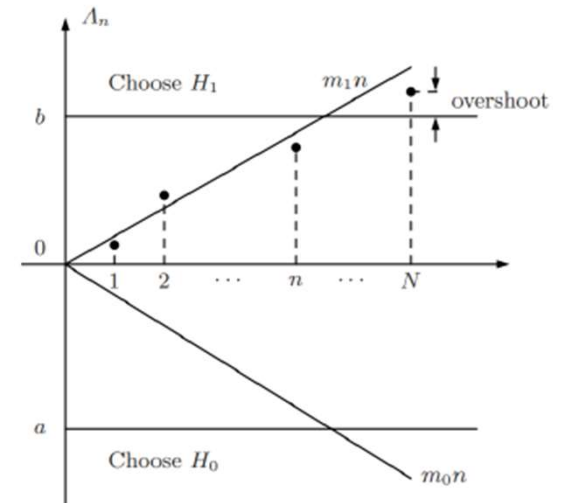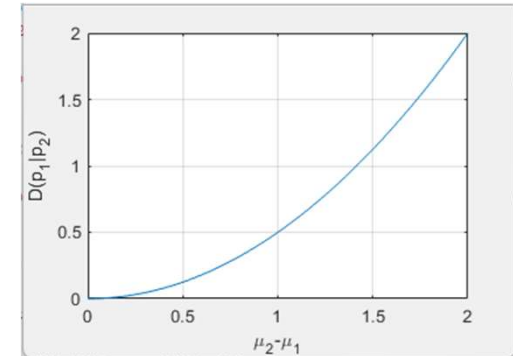  – **Subject to:** Constraint on false positive rate

$$\mathbb{E}^{(\infty)}[d^\kappa] \geq \tau$$



Kim et al, "Human-in-the-loop Cyber Intrusion Detection Using Active Learning" IEEE TIFS, 2024

# Background: Sequential Hypothesis Testing

- Generalized likelihood ratio test
  - Composite hypothesis $\mathcal{H}_1 = \{h_1, \ldots, h_H\}$
  - Detection rule $\eta_t = \begin{cases} \mathcal{H}_1 \; if \; \dfrac{\underset{h \in \mathcal{H}_1}{\max} P(Y|h)}{P(Y|h_0)} > \Theta \\ \quad\quad otherwise \end{cases}$

- Asymptotic behavior
  - Risk
    - $R_h \triangleq \underset{h' \neq h}{\max} P_{h'}(\eta_t = h)$
  - Expected detection time
    - $E\left[t_d^h\right] \geq \dfrac{-log R_h}{D(p_h||p_{h'})}(1 + o(1))$

B. C. Levy, Principles of signal detection and parameter estimation. Springer, 2008.

# Defender's Problem

- **Defender objective:** Minimize mean time to detection

$$\kappa^* = \arg\min_{\kappa \in \mathcal{K}} \sup_{t_{1 \to 2} > 0} \mathbb{E}^{(t_{1 \to 2})}[d^\kappa - t_{1 \to 2}]$$

  - **Subject to:** Constraint on false positive rate
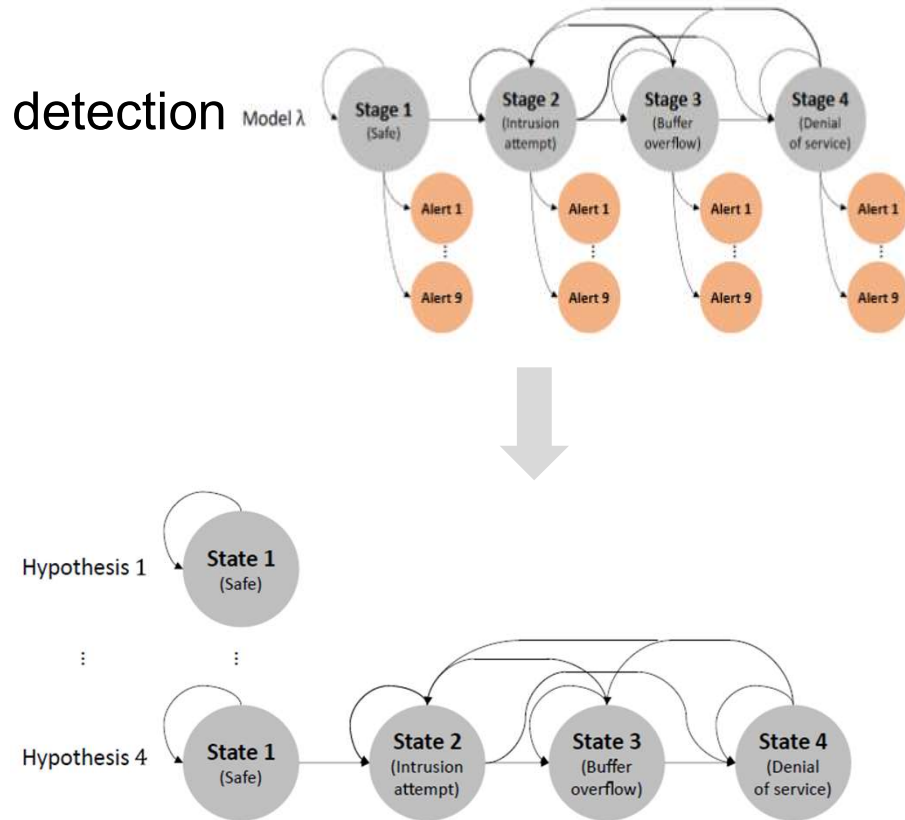
$$\mathbb{E}^{(\infty)}[d^\kappa] \geq \tau.$$

- **Generating Alternative Hypotheses**

  - Most likely hypothesis at time $t$

$$\hat{h} = argmax_{h \in \mathcal{H}} P_h(Y_{1:t} | \mathcal{F}_t, v_t^\kappa)$$

  - Likelihood ratio

> $S_t^\kappa = \dfrac{P_{\hat{h}}(Y_t | \mathcal{F}_t, v_t^\kappa)}{P_1(Y_t | \mathcal{F}_t, v_t^\kappa)}$

# Active Learning for Quickest Detection

- **Optimal detection rule without active learning**

  - Generalized likelihood ratio test

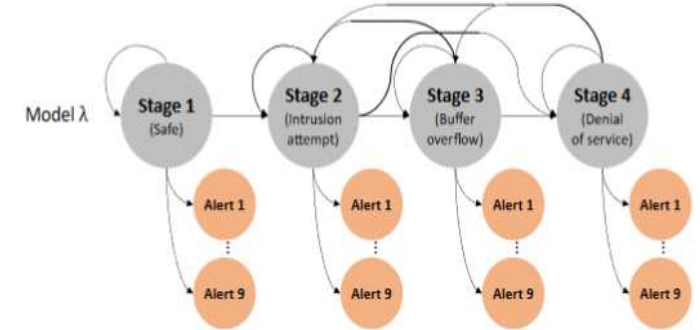- **Two candidate policies**

  - Max-ratio policy

    → Set of alerts that maximizes the expected probability ratio

$$\mathcal{V}_t^{MR} = \operatorname*{arg\,max}_{v_t \subseteq Y_{1:t}^+, |v_t| \leq B} \left| \mathbb{E}\left[ \frac{p_{\hat{h}}(Y_{1:t} = y_{1:t} | \mathcal{F}_t = f_t, \mathcal{V}_t = v_t)}{p_1(Y_{1:t} = y_{1:t} | \mathcal{F}_t = f_t, \mathcal{V}_t = v_t)} \right] \right|$$
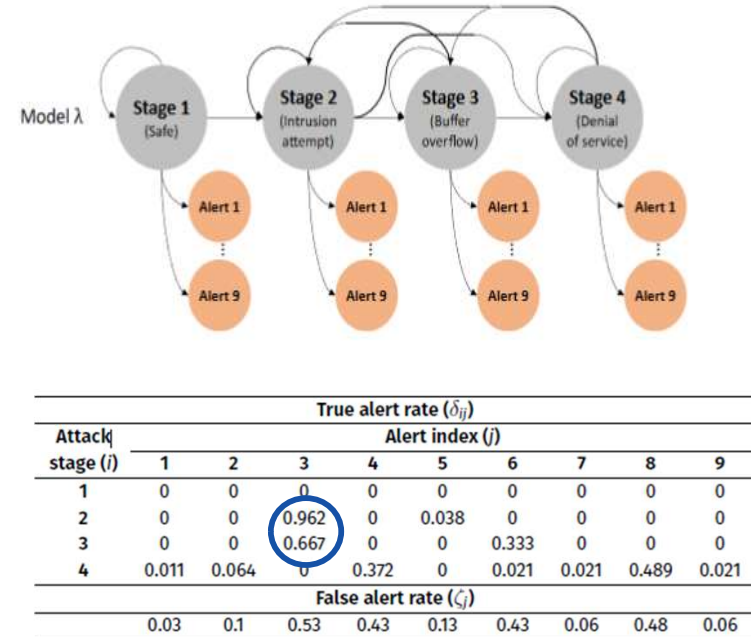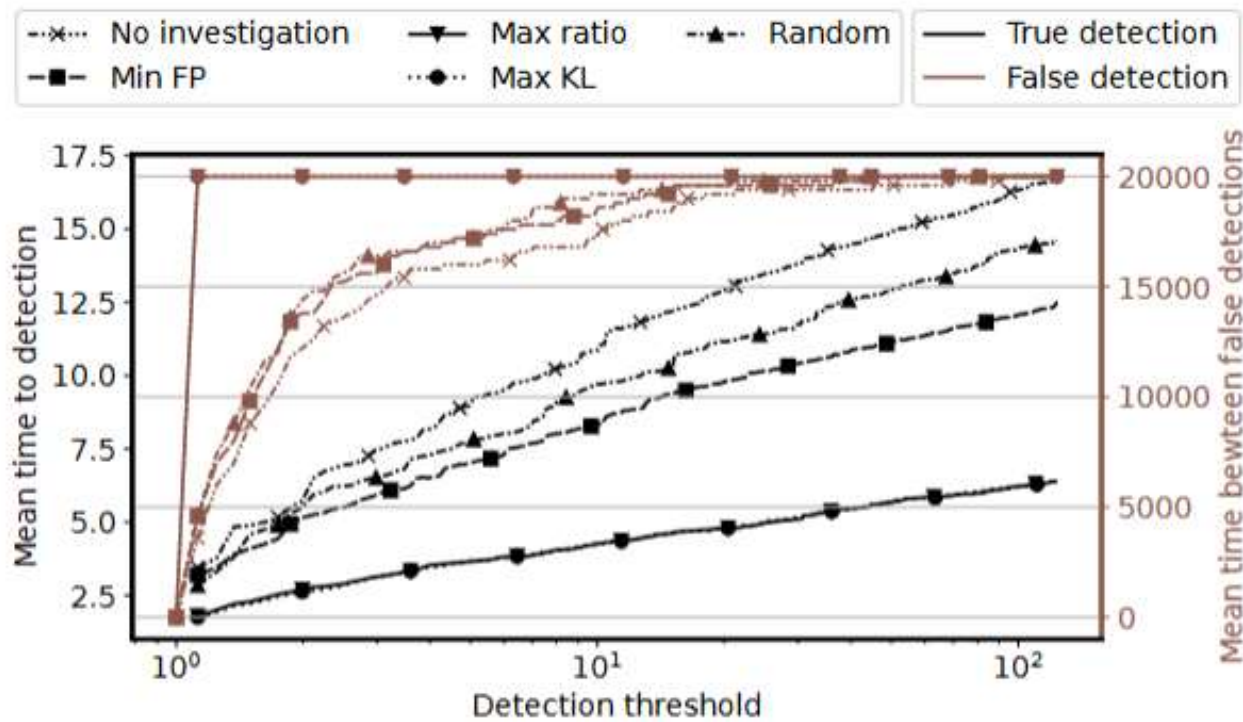
  - Max KL Divergence

    → Set of alerts that maximize the KL divergence of the distribution of observed alerts after investigation
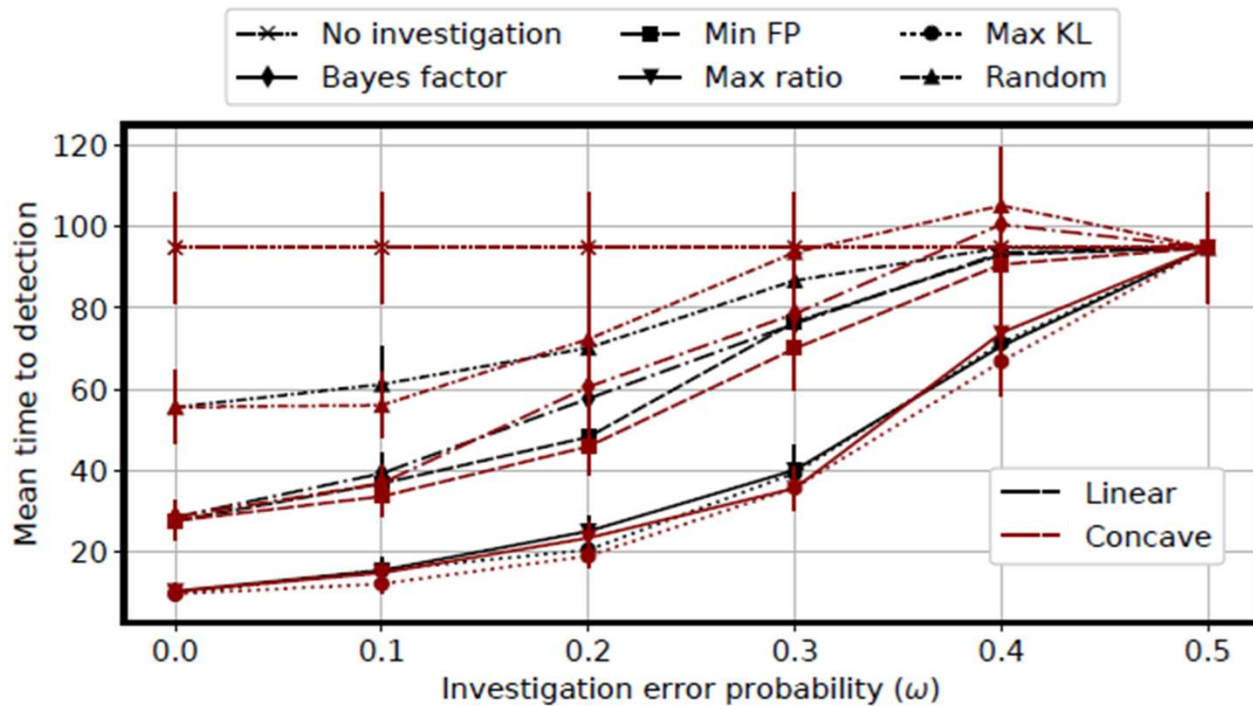
$$\mathcal{V}_t^{MKL} = \operatorname*{arg\,max}_{v_t \subseteq Y_{1:t}^+, |v_t| \leq B} \mathbb{E}\Big[ \sum_{t'=1}^{t} D\big(\mathbb{P}_{\hat{h}}(Y_{t'} = y_{t'} | \mathcal{F}_t = f_t, \mathcal{V}_t = v_t) \,\|\, \mathbb{P}_1(Y_{t'} = y_{t'} | \mathcal{F}_t = f_t, \mathcal{V}_t = v_t))\Big]$$
$$- \sum_{t'=1}^{t} D(\mathbb{P}_{\hat{h}}(Y_{t'} = y_{t'} | \mathcal{F}_t = f_t) \| \mathbb{P}_1(Y_{t'} = y_{t'} | \mathcal{F}_t = f_t))$$
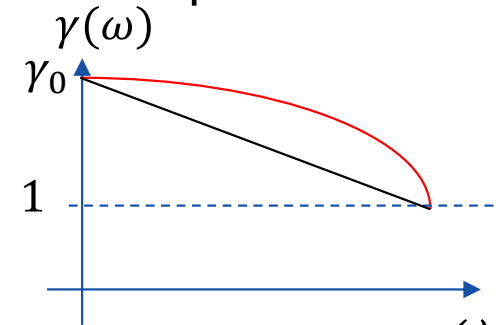
Kim et al, "Human-in-the-loop Cyber Intrusion Detection Using Active Learning" IEEE TIFS, 2024

# Detection Performance
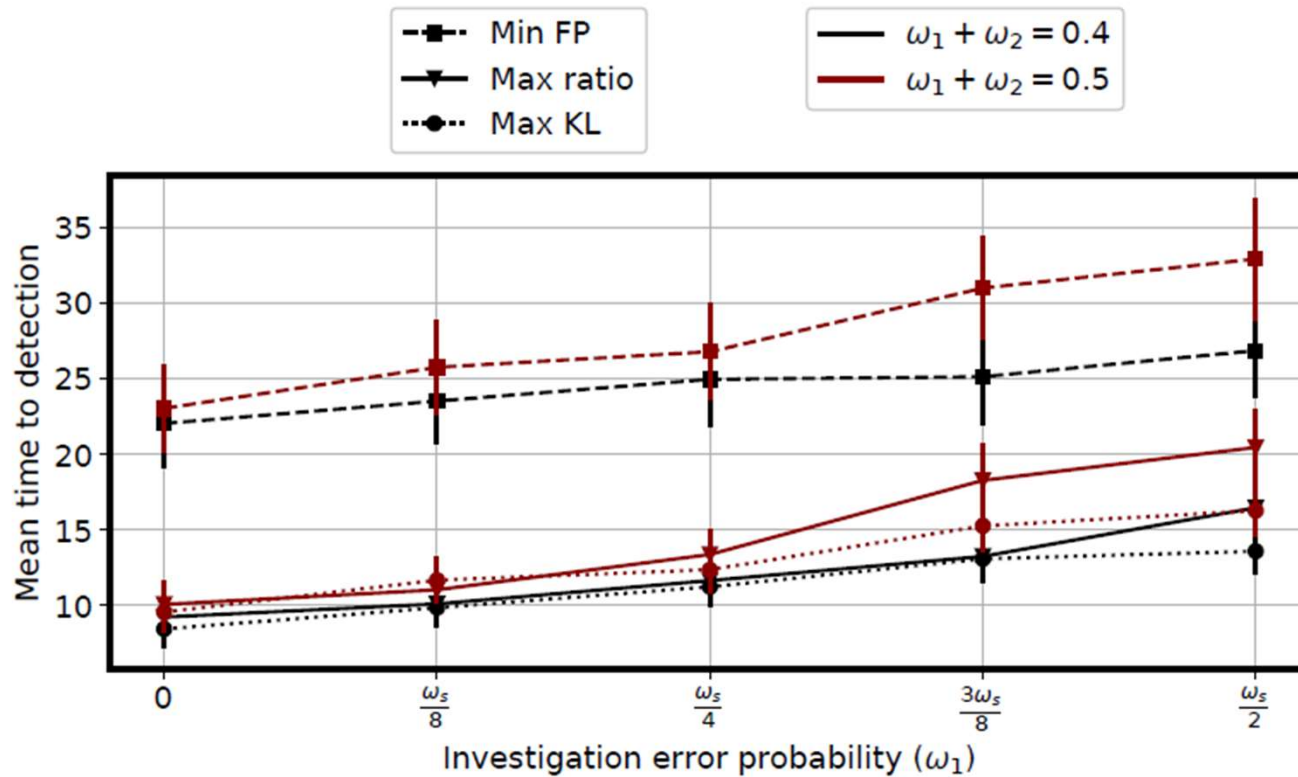
# Impact of the Human Model
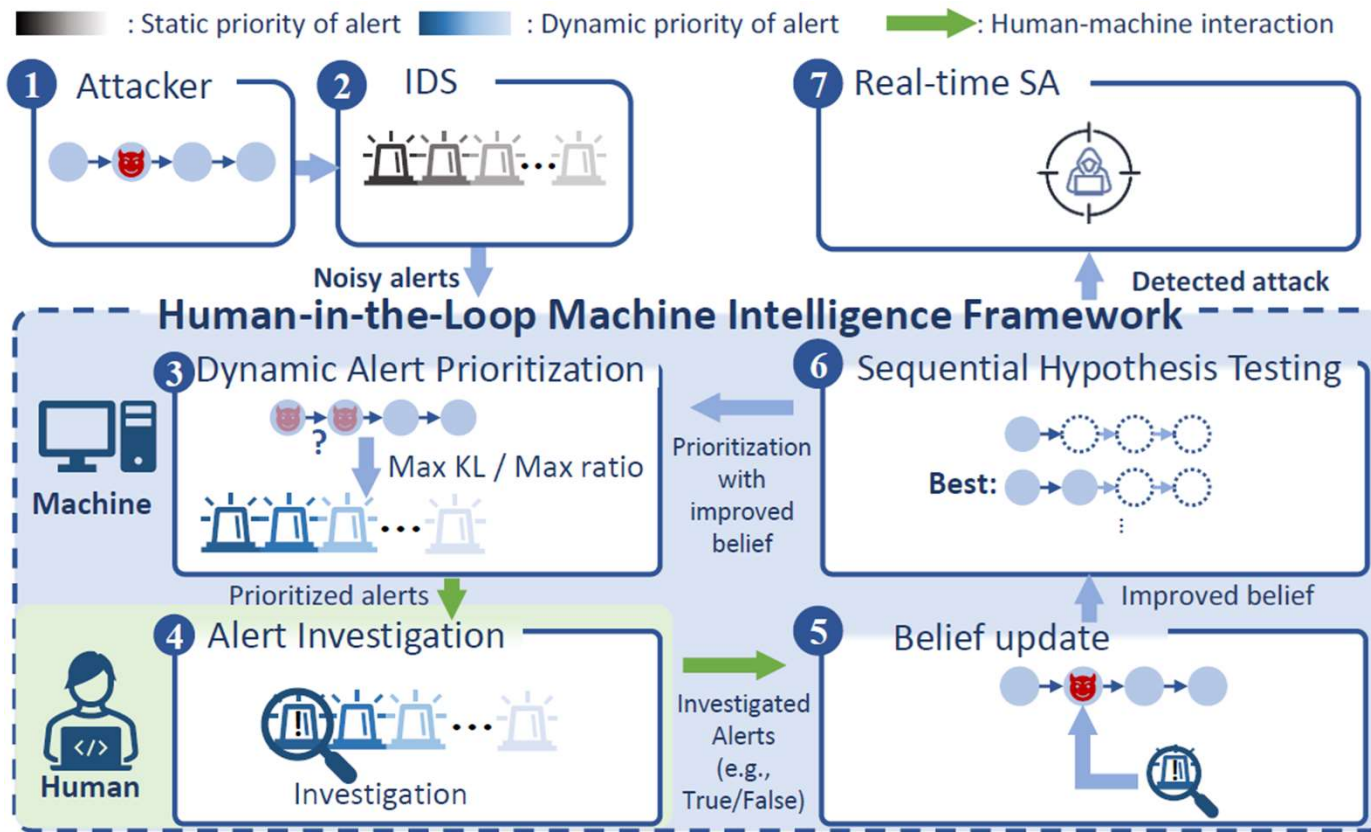


- Concave confidence function superior

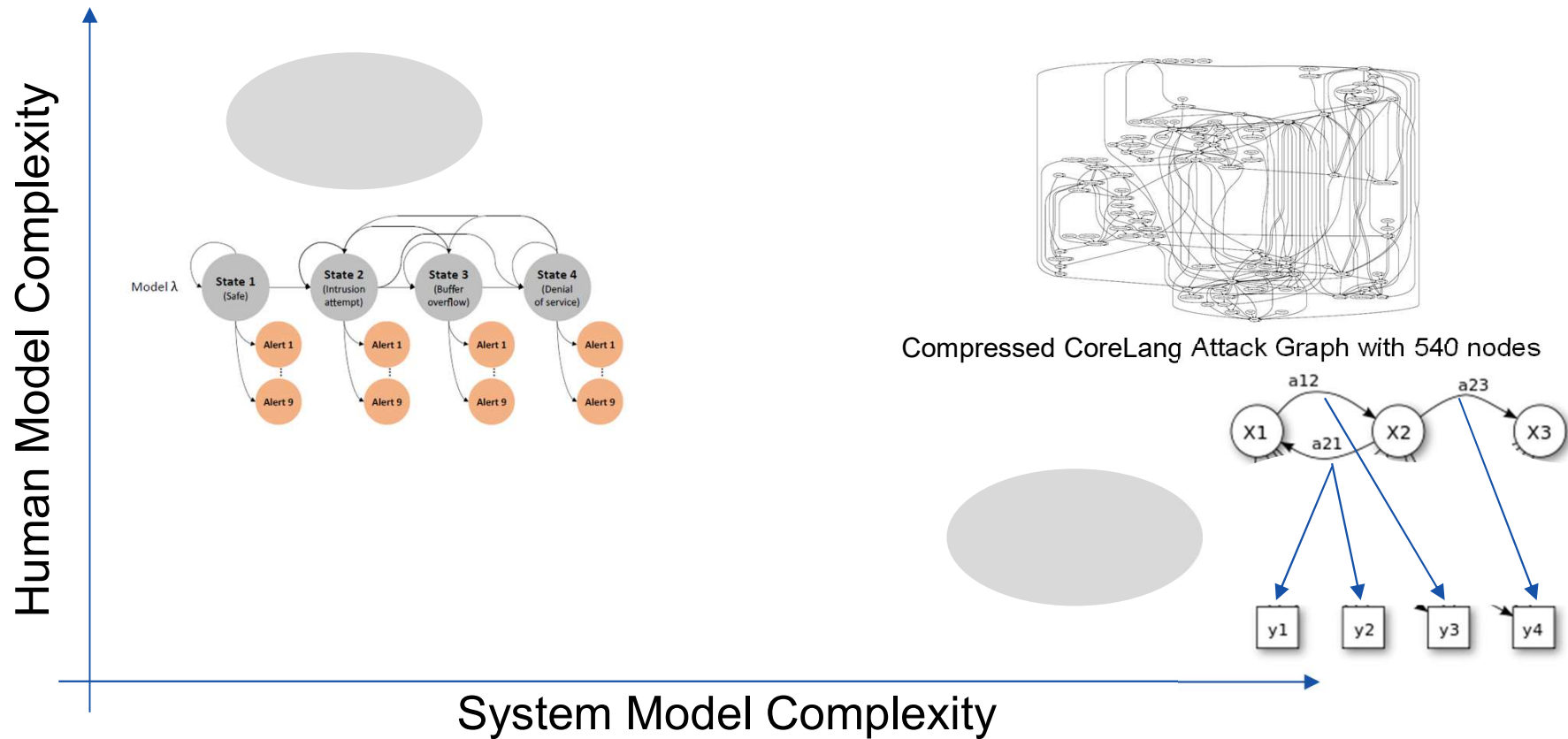$$\gamma(\omega)$$

- Max KL performs best

# Expertise is Important



- 2 experts with potentially varying expertise
- Heterogenous expertise is preferrable

Kim et al, "Human-in-the-loop Cyber Intrusion Detection Using Active Learning" IEEE TIFS, 2024

# Human-in-the-Loop AI Framework Revisited

# Framework Design Space



Human Model Complexity *(vertical axis)*

System Model Complexity *(horizontal axis)*

Compressed CoreLang Attack Graph with 540 nodes

Katsikeas et al. "An attack simulation language for the IT domain," in *Proc. of Int. Workshop on Graphical Models for Security*, pp. 67–86, 2020.
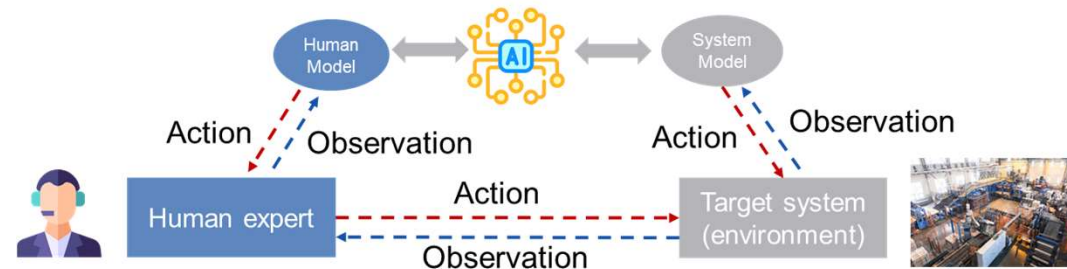
# Conclusion

- **Human-in-the-loop-AI for cyber resilience**

  - Efficient use of human resources and ML

  - Human skills and behavior vs. system model complexity

  - Improved accuracy and lower time to detection

- **Many open questions**

  - How to model human behaviour

    > *Trust, psychological aspects*

    > Affects the design of AI algorithms

  - How to apply the concept to CPS

  - Vulnerability to an adaptive adversary in a game theoretical framework

  - Integration with threat hunting

  - Semi-autonomous incident response

# Boosting Cyber Resilience with Human-in-the-loop AI

György Dán

IEEE CNS 2024 Workshop on Cyber Resilience