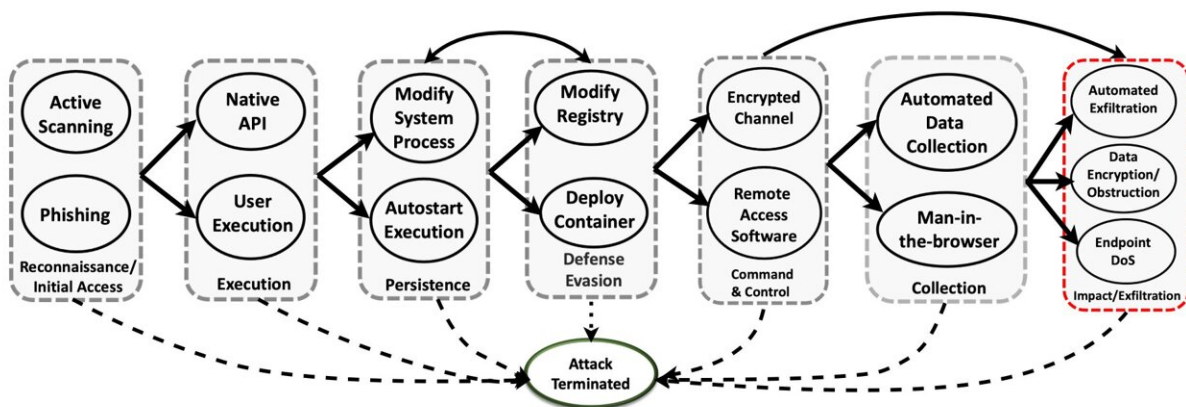


Cybersecurity defenders are expanding their AI toolbox

February 16 2023



Multi-stage attack propagation represented with MITRE ATT&CK Tactics and Techniques. (Note: A directed edge between an attack tactic and technique specifies that the attacker may try to implement that technique next after achieving the objective of the attack tactic. Bidirectional arrow represents that Defense Evasion can come before Persistence.). Credit: *arXiv* (2023). DOI: 10.48550/arxiv.2302.01595

Scientists have taken a key step toward harnessing a form of artificial intelligence known as deep reinforcement learning, or DRL, to protect computer networks.

When faced with sophisticated cyberattacks in a rigorous simulation setting, [deep reinforcement learning](#) was effective at stopping

adversaries from reaching their goals up to 95 percent of the time. The outcome offers promise for a role for autonomous AI in proactive cyber defense.

Scientists from the Department of Energy's Pacific Northwest National Laboratory documented their findings in a research paper and presented their work Feb. 14 at a workshop on AI for Cybersecurity during the annual meeting of the Association for the Advancement of Artificial Intelligence in Washington, D.C.

The starting point was the development of a simulation environment to test multistage attack scenarios involving distinct types of adversaries. Creation of such a dynamic attack-defense simulation environment for experimentation itself is a win. The environment offers researchers a way to compare the effectiveness of different AI-based defensive methods under controlled test settings.

Such tools are essential for evaluating the performance of deep reinforcement learning algorithms. The method is emerging as a powerful decision-support tool for cybersecurity experts—a defense agent with the ability to learn, adapt to quickly changing circumstances, and make decisions autonomously. While other forms of AI are standard to detect intrusions or filter spam messages, deep reinforcement learning expands defenders' abilities to orchestrate sequential decision-making plans in their daily face-off with adversaries.

Deep reinforcement learning offers smarter cybersecurity, the ability to detect changes in the cyber landscape earlier, and the opportunity to take preemptive steps to scuttle a cyberattack.

DRL: Decisions in a broad attack space

"An effective AI agent for cybersecurity needs to sense, perceive, act

and adapt, based on the information it can gather and on the results of decisions that it enacts," said Samrat Chatterjee, a data scientist who presented the team's work. "Deep reinforcement learning holds great potential in this space, where the number of system states and action choices can be large."

DRL, which combines reinforcement learning and deep learning, is especially adept in situations where a series of decisions in a complex environment need to be made. Good decisions leading to desirable results are reinforced with a positive reward (expressed as a numeric value); bad choices leading to undesirable outcomes are discouraged via a negative cost.

It's similar to how people learn many tasks. A child who does their chores might receive [positive reinforcement](#) with a desired playdate; a child who doesn't do their work gets negative reinforcement, like the takeaway of a digital device.

"It's the same concept in reinforcement learning," Chatterjee said. "The agent can choose from a set of actions. With each action comes feedback, good or bad, that becomes part of its memory. There's an interplay between exploring new opportunities and exploiting past experiences. The goal is to create an agent that learns to make [good decisions](#)."

Open AI Gym and MITRE ATT&CK

The team used an open-source software toolkit known as Open AI Gym as a basis to create a custom and controlled simulation environment to evaluate the strengths and weaknesses of four deep reinforcement learning algorithms.

The team used the MITRE ATT&CK framework, developed by MITRE

Corp., and incorporated seven tactics and 15 techniques deployed by three distinct adversaries. Defenders were equipped with 23 mitigation actions to try to halt or prevent the progression of an attack.

Stages of the attack included tactics of reconnaissance, execution, persistence, defense evasion, command and control, collection and exfiltration (when data is transferred out of the system). An attack was recorded as a win for the adversary if they successfully reached the final exfiltration stage.

"Our algorithms operate in a competitive environment—a contest with an adversary intent on breaching the system," said Chatterjee. "It's a multistage attack, where the adversary can pursue multiple attack paths that can change over time as they try to go from reconnaissance to exploitation. Our challenge is to show how defenses based on deep reinforcement learning can stop such an attack."

DQN outpaces other approaches

The team trained defensive agents based on four deep [reinforcement learning](#) algorithms: DQN (Deep Q-Network) and three variations of what's known as the actor-critic approach. The agents were trained with simulated data about cyberattacks, then tested against attacks that they had not observed in training.

DQN performed the best.

- Least sophisticated attacks (based on varying levels of adversary skill and persistence): DQN stopped 79 percent of attacks midway through attack stages and 93 percent by the final stage.
- Moderately sophisticated attacks: DQN stopped 82 percent of attacks midway and 95 percent by the final stage.
- Most sophisticated attacks: DQN stopped 57 percent of attacks

midway and 84 percent by the final stage—far higher than the other three algorithms.

"Our goal is to create an autonomous defense agent that can learn the most likely next step of an adversary, plan for it, and then respond in the best way to protect the system," Chatterjee said.

Despite the progress, no one is ready to entrust cyber defense entirely up to an AI system. Instead, a DRL-based cybersecurity system would need to work in concert with humans, said co-author Arnab Bhattacharya, formerly of PNNL.

"AI can be good at defending against a specific strategy but isn't as good at understanding all the approaches an adversary might take," Bhattacharya said. "We are nowhere near the stage where AI can replace human cyber analysts. Human feedback and guidance are important."

The research is published on the *arXiv* preprint server.

More information: Ashutosh Dutta et al, Deep Reinforcement Learning for Cyber System Defense under Dynamic Adversarial Uncertainties, *arXiv* (2023). [DOI: 10.48550/arxiv.2302.01595](https://doi.org/10.48550/arxiv.2302.01595)

Provided by Pacific Northwest National Laboratory

Citation: Cybersecurity defenders are expanding their AI toolbox (2023, February 16) retrieved 26 April 2024 from <https://techxplore.com/news/2023-02-cybersecurity-defenders-ai-toolbox.html>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.