

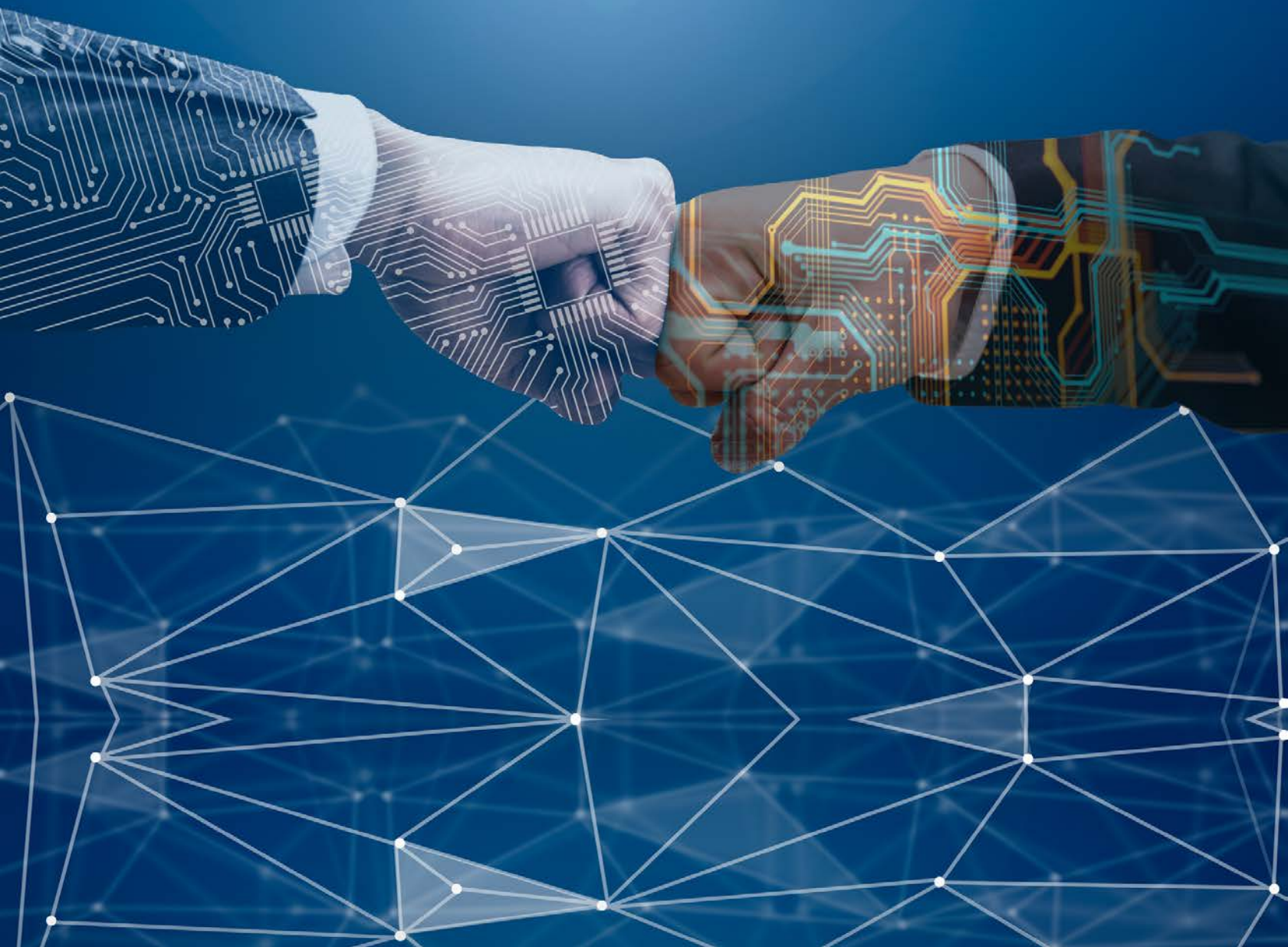
HOT OFF THE PRESS



When AI Fights Back:

How AI Is Rewriting Cyber Defense

By Sandy Winnefeld



When AI Fights Back: How AI Is Rewriting Cyber Defense

By Sandy Winnefeld

The Moment the Trap Snapped

The first time I watched an AI attack collapse on itself, I realized the defenders hadn't just blocked it, they'd tricked it. The network lied so convincingly that the machine couldn't tell what was real.

A few minutes earlier, I'd been observing a simulated spear-phishing attack that crafted a convincing email with help from a large language model. I watched an unwitting employee click the link and open the door to the company's internal network. From there, the plan was simple: map, move, and capture. But the system had other ideas.

Nothing behaved as expected. The files looked real but weren't. The servers responded, then vanished. Even the credentials I harvested dissolved upon use. Within seconds, I was locked out, unsure what had been real and what had been bait.

The network hadn't beaten the attacking team with speed or firepower. It had fooled them.

That's when I understood: the defenders had deployed AI-powered deception, a counterintelligence system built to mislead intruders rather than chase them.

The Age of AI Offense

Artificial intelligence has given attackers new capabilities that humans alone could never scale. Phishing kits now write their own lures. Ransomware builders auto-mutate code to escape detection. Deepfake voices impersonate executives with unnerving accuracy.

What once took days of reconnaissance can now happen in seconds, driven by algorithms that learn, adjust, and strike before defenders can even analyze what hit them. Traditional prevention tools—firewalls, antivirus, endpoint agents—still matter, but they rely on what they have already seen.

AI has made novelty a weapon. Each attack can be unique, synthetic, and perfectly timed. Detection systems built on pattern recognition are facing something that constantly invents new patterns. That's a fight they can't win on speed alone.

When AI Fights Back: How AI Is Rewriting Cyber Defense

By Sandy Winnefeld

Turning Deception Into Defense

To fight an intelligent opponent, defenders are beginning to think more like tricksters than guards. Deception technology plants digital illusions—decoy servers, fake credentials, synthetic data—that look authentic enough to lure both human hackers and AI reconnaissance systems. Once an attacker touches one of these false assets, the defense gains precise telemetry on tactics, timing, and intent.

It's an old idea updated for a machine age. Militaries have used deception for centuries, but applying it inside live networks requires realism at machine scale. That realism now comes from AI itself. Honeypots and the like used to be clunky, easy to detect, and not producible at scale (because humans had to do it). That has all changed.

Modern deception systems can automatically create thousands of believable targets, each subtly different and able to adapt when probed. The result is a living labyrinth where every step could be a trap.

When the Algorithms Collide

AI has become both attacker and defender, locked in a duel of misdirection. Attack bots crawl systems, parsing responses and searching for weaknesses. Deception AIs answer back with convincing noise—files that look valuable, systems that behave normally, and networks that shift faster than reconnaissance scripts can adapt.

The attacker's confidence erodes as the data becomes unreliable. In effect, the defender breaks the machine's model of reality. This contest plays out at machine speed, but its outcome depends on something far older: uncertainty.

This is more than cybersecurity. It's a philosophical shift. For decades, defense meant control and certainty. Now, survival may depend on ambiguity. The goal is no longer to seal every door but to make the intruder doubt which doors exist.

When AI Fights Back: How AI Is Rewriting Cyber Defense



By Sandy Winnefeld

The Human Consequence

The collision between offensive and defensive AI raises deeper questions. What happens when synthetic intelligence becomes the primary actor in digital conflict? If machines can deceive one another convincingly, how do humans maintain situational awareness?

Already, the same generative systems that protect networks are being trained to fabricate documents, impersonate users, and wage information warfare at scale. The boundary between attack and counterattack is blurring.

Every defensive deception risks teaching offensive models new tricks. Every dataset used to train an AI system becomes a potential manipulation target itself. What began as a technical contest of code is turning into a test of trust.

A New Philosophy of Defense

For organizations, the lesson is clear: prevention alone cannot keep up. AI-driven threats demand AI-driven defense, but one that adds an element human attackers never mastered—psychological uncertainty.

Deception forces even autonomous systems to hesitate, to verify, to lose the speed advantage that defines them. The future of cybersecurity may not be walls and sensors but mirrors and mazes. Truth, even in data, is becoming negotiable, and defense now depends on who controls the illusion.



Admiral (Ret.) James "Sandy" Winnefeld retired as the 9th Vice Chairman of the Joint Chiefs of Staff, and is on Acalvio's Federal Advisory Board. His interests and investments include the field of AI-enabled active defense.

Acalvio is the leader in autonomous cyber deception, defending against APTs, insider threats, and ransomware. Its AI-powered Preemptive Cybersecurity Platform, protected by 25 patents, delivers threat detection across IT, OT, and cloud environments and advances Identity Threat Detection and Response (ITDR) with Honeytokens-driven Zero Trust security. Based in Silicon Valley, Acalvio serves midsize to Fortune 500 companies and government agencies, offering flexible deployment from the Cloud, on-premises, or through managed service providers. www.acalvio.com