

AI-ENABLED THREAT
ACTORS AND PURPLE
TEAMING: A STRATEGIC
EXECUTIVE FRAMEWORK

PENTRA SECURITY WHITE PAPER

pentrasecurity.com

About Us



We are the architects of world-class security programs, specializing in purple teaming, offensive security, and incident response. Our proven methodologies transform fragmented security operations into unified, measurable defense capabilities that withstand real-world attacks.

What We Do

We provide comprehensive purple teaming programs that combine offensive and defensive capabilities into cohesive security operations. Our services span the entire security lifecycle, from initial attack simulation and vulnerability assessment to incident response and continuous improvement. We provide clear, repeatable methodologies that transform security teams from siloed operations into synchronized defense forces capable of detecting and defending against advanced threats at machine speed.

Who We Are

We are seasoned security professionals who have led offensive and defensive operations at the highest levels of industry and government. Our team combines decades of red team experience with blue team operational excellence, united by a shared commitment to collaborative security that delivers measurable results. Each consultant has hands-on experience in building and executing world-class security programs from the ground up.



Our goal is to revolutionize cybersecurity through collaborative purple teaming that proves defense capabilities against real-world attacks. We empower organizations to measure, validate and continuously improve their security posture by creating repeatable frameworks that lead to quantifiable risk mitigation. By combining red and blue perspectives, we create security programs that anticipate threats instead of reacting to breaches.



We envision a future where every organization operates with the assurance that its security has been put to the test through continuous validation by a purple team. Our vision goes beyond traditional penetration testing and establishes purple teaming as the cornerstone of modern security operations. We are building the next generation of security programs where offensive innovation and defensive excellence come together to create impenetrable digital fortresses.

Executive Summary



The cybersecurity landscape has changed fundamentally. Al-powered threat actors are now responsible for 67.4% of all phishing attacks, with attack volumes increasing by over 4,000% since the introduction of ChatGPT. Organizations are facing unprecedented challenges as criminals weaponize artificial intelligence to scale attacks, lower barriers to entry and improve sophistication at machine speed.

By the Numbers

AI-Enabled **67.4%**

Al-enabled threat actors now drive 67.4% of all phishing attacks

Outperform 55%

Al-powered phishing campaigns outperform elite cybercriminals by 55% Voice Phishing 442%

Phishing attacks increased 442% between the first and second half of 2024.

DarkGPT **\$2.85m**

FraudGPT command \$200-\$1,700 annual subscriptions with over 3,000 confirmed sales \$4.48m

The average cost of Al-enhanced phishing breaches now reaches \$4.88 million

Purple Teaming 40%

Mature purple team programs report 40% reductions in mean time to detect

Current security frameworks, while useful, have critical gaps that leave organizations vulnerable. MITRE ATLAS provides tactical information but lacks guidance for operational implementation. The NIST AI RMF provides strategic direction, but requires significant organizational maturity that most companies lack. This research shows that there is an urgent need for a new methodology that combines strategic frameworks with operational security while addressing AI-specific threats that traditional cybersecurity cannot adequately address.

Purple Teaming, the collaborative integration of offensive Red Team capabilities with defensive Blue Team expertise, has proven to be the most effective approach to combat Al-powered threats. Organizations implementing mature Purple Team programs achieve a 3:1 to 5:1 return on their security investment, with a 40% improvement in threat detection rates and significantly reduced incident response times. However, existing purple teaming methods need to be fundamentally adapted to take into account Al-specific attack vectors and defense requirements.

The convergence of rapidly evolving AI threats, regulatory pressures from the EU AI Act and emerging compliance frameworks, and the proven effectiveness of Purple Teaming create a strategic opportunity. This research proposes a comprehensive new methodology that goes beyond current frameworks to provide leaders with practical, measurable and customizable AI safety capabilities.



Current state of Al-enabled threats demands immediate action

The threat landscape has evolved with alarming speed and sophistication. Research from leading cybersecurity firms shows that nation-state actors in five countries are using Al operationally for cyber operations, while criminal organizations are using Al tools to generate sophisticated malware in minutes that previously required expert knowledge.

AI ENABLED ATTACKS

Deepfake \$25m

A single deepfake CEO fraud resulted in \$25 million in losses, impersonating senior execs.

Al-Enhanced

\$866,255

Generated at least \$866,255 from just 10 of 64+ compromised organizations.

Companies

320

Al-enhanced attacks, infiltrated over 320 companies

The scale of the use of AI by threat actors is staggering. Voice phishing attacks increased by 442% between the first and second half of 2024, driven by AI-generated synthetic voices that perfectly mimic executives and trusted contacts. The average cost of AI-powered phishing attacks now stands at \$4.88 million, a 10% increase over traditional attacks. Of particular concern is that AI-powered phishing campaigns outperform the effectiveness of elite human cybercriminals by 55%.

Real-world incidents show the immediate impact on business. A single deepfake CEO scam resulted in \$25 million in losses when multiple executives simultaneously impersonated them in a video conference. North Korean IT operatives used Al-powered attacks to infiltrate over 320 companies, generating at least \$866,255 in just 10 of more than 64 compromised organizations. These attacks were successful because they used Al to accelerate traditional attack techniques while introducing entirely new threat vectors that existing security controls cannot adequately address.

Criminal AI tools have rapidly professionalized. Dark web platforms like FraudGPT charge annual subscriptions worth \$200 to \$1,700 and have over 3,000 confirmed sales. These platforms offer undetectable malware generation, automated credential harvesting and social engineering content creation without ethical guardrails. The democratization of advanced attack capabilities means that even inexperienced criminals can now execute sophisticated campaigns that previously required years of experience.



Purple teaming emerges as the most effective defense strategy

Traditional red team penetration testing and blue team defenses, while independently valuable, cannot adequately respond to Al-powered threats. The speed, scale and sophistication of Al-powered attacks require collaborative security measures that combine offensive simulations with real-time defensive improvements.

Purple Teaming has demonstrated measurable business value across multiple industries. Companies implementing mature Purple Team programs report a 40 percent reduction in average time to threat detection, with some achieving as much as 99 percent improvements in both detection and response time. Financial benefits include an average \$3 million reduction in data breach costs for organizations with fully implemented security AI capabilities.

PURPLE TEAMING

Purple Teaming 99%

Achieving 99% improvements in both detection and response times

Breach Reduction

Financial benefits include an average \$3 million reduction in data breach costs

Post-Pentesting 40%

On average Purple Teaming can quickly identify and close gaps from a Pentest

The Mitek Systems case study illustrates the potential of Purple Teaming for AI security. The award-winning purple team focuses exclusively on AI-specific attacks, using cutting-edge techniques such as 3D masks, GAN-generated deepfakes and document spoofing to test their fraud detection systems. This approach has positioned their AI security capabilities among the industry leaders while providing continuous validation of their defensive improvements.

However, traditional purple teaming methods need to be significantly adapted for AI threats. Existing frameworks such as SANS SEC599 and SCYTHE's Purple Team Exercise Framework address conventional cybersecurity, but lack comprehensive guidance for AI-specific attack vectors such as adversarial examples, data poisoning, prompt injection and model extraction attacks. Organizations need specialized approaches that address both conventional cybersecurity and the vulnerabilities of AI/ML systems.



Current frameworks fall short of addressing Al security needs

A comprehensive analysis of existing frameworks reveals significant implementation gaps that represent an organizational weakness. While MITRE ATLAS provides excellent tactical threat intelligence, it lacks practical guidance for translating threats into operational security programs. Organizations struggle with risk prioritization, control mapping and measuring the effectiveness of the framework. The framework's highly technical focus creates barriers to executive engagement and strategic decision making.

ATTACK VOLUME

23%

Only 23% of organizations successfully operationalize MITRE ATLAS Significant Gaps 84%

84% of CISOs report significant gaps between NIST AI RMF strategic guidance and their operational security capabilities Al Security 2.7x

Organizations using multiple AI security frameworks spend 2.7x more on compliance overhead

The NIST AI Risk Management Framework provides strategic direction, but requires sophisticated risk management capabilities that most organizations have not yet developed. The framework requires significant investment in governance infrastructure and provides limited technical specificity for implementing security controls. Organizations report difficulties in quantifying AI risks and measuring the effectiveness of the framework.

The sector-specific frameworks remain fragmented and inconsistent. Financial services regulations include multiple authorities with different requirements. Healthcare frameworks expand on existing HIPAA compliance without addressing AI-specific risks. Critical infrastructure guidelines emphasize sector-specific assessments but provide limited guidance on operational security.

Vendor-specific frameworks create additional complications. Cloud providers offer comprehensive security features but lock companies into specific platforms. Commercial Al security tools offer point solutions without integrated approaches. The result is a complex ecosystem with overlapping requirements without consistent implementation guidelines.

The analysis identifies clear opportunities for a new methodology that closes these gaps through practical implementation bridges, automated assessment capabilities and leadership integration. Organizations need frameworks that translate strategic guidance into operational certainty while ensuring measurable business value and regulatory compliance.



Al supply chain vulnerabilities create unprecedented risk exposure

The AI ecosystem comprises data sets, models and software components that create complex dependencies and unique attack surfaces. Research shows that poisoning as little as 0.001% of training data can lead to compromised models, while the rise of malicious AI models circulating in open source repositories has increased by more than 1,300. A simple example of this would be an attacker adding a set of pixels to images of dogs that when scanned by an object detection model that tricks it into thinking that an enemy tank is a Doberman.

AI VULNERABILITIES

Poisoning 0.001%

Poisoning just 0.001% of training data can result in compromised models

Exposure Points 1300%

1,300% increase in malicious AI models circulating repositories creates exposure points

Vulnerable Al 92%

92% of organizations use at least one vulnerable Al component

Attacks on the AI supply chain can have a cascading effect on entire systems. Model namespace reuse attacks on platforms such as Hugging Face allow attackers to use malicious substitutes for legitimate AI models. Third-party dependencies in AI frameworks create vulnerabilities that can compromise entire ML pipelines. Critical dependency on open source datasets such as LAION 5B creates widespread vulnerabilities that cannot be adequately monitored with traditional security controls.

Cloud vendors have developed comprehensive AI security frameworks, but implementation requires deep technical expertise and significant investment in resources. AWS's Generative AI Security Scoping Matrix, Microsoft's AI Foundry Security Baseline and Google's Secure AI Framework provide valuable starting points, but require an organizational maturity that many companies lack.

MLOps security is a critical gap in most companies. Securing Al workflows requires data lineage tracking, model version control, secure development environments and runtime monitoring capabilities. Container security, Kubernetes management and handling secrets for Al workloads require specialized expertise that combines traditional DevSecOps with Al/ML domain knowledge.



Regulatory compliance creates urgent implementation requirements

The regulatory landscape requires immediate executive attention and strategic planning. The phased implementation of the EU AI law by 2027 creates mandatory requirements for companies using AI systems and provides for significant penalties of up to 7% of global turnover for non-compliance. High-risk AI systems must demonstrate an "appropriate level of accuracy, robustness and cybersecurity", with mandatory reporting within two days for serious security incidents.

COMPLIANCE

Al Regulations 73%

73% of multinational corporations face conflicting Al regulations

Al Incidents
48 hours

Organizations have only 48 hours to report serious Al security incidents Non-Compliance 4.7%

Non-compliance penalties average 4.7% of global annual revenue

Following the Trump administration's rescission of Executive Order 14110, regulatory approaches in the US have shifted significantly, but NIST continues to develop voluntary frameworks and industry-specific regulations remain active. Financial services are subject to enhanced due diligence requirements, healthcare organizations must comply with Al-specific HIPAA extensions, and critical infrastructure operators are subject to enhanced cybersecurity requirements for Al-dependent systems.

International coordination through the G7 Hiroshima process provides a voluntary framework, but the development of technical standards is still spread across multiple organizations. Global companies are faced with complex compliance requirements that span multiple jurisdictions with different technical standards and enforcement approaches.

The analysis shows that compliance cannot be achieved through stand-alone solutions, but requires integrated governance systems that span the entire Al lifecycle. Organizations need frameworks that provide continuous compliance monitoring, automated reporting capabilities and clear accountability structures that meet both strategic governance requirements and operational security requirements.



Novel approaches point toward next-generation methodologies

Research from leading universities, innovative start-ups and government programs reveals cutting-edge defense approaches that go well beyond current frameworks. The MAESTRO framework represents the most significant advance. It was developed specifically for agent-based AI and has a seven-layer architecture that covers vulnerabilities from base models to agent ecosystems.

AI POWERED

Al-Powered 94%

Al-powered defense systems reduce threat detection time by 94% Zero-Trust Al 78%

Organizations
implementing Zero
Trust Al architectures
experience 78% fewer
successful attacks

Al Detection 77%

Al systems achieving 77% vulnerability detection rates with 61% automated patching success

DARPA's AI Cyber Challenge has shown that AI systems can detect 77% of vulnerabilities and automatically patch 61% of vulnerabilities at a cost of \$152 per task compared to thousands of dollars for traditional methods. These results prove that AI-powered defense systems can achieve human expert-level effectiveness while operating at machine speed and scale.

Advanced behavioral analytics using Indicators of Behavior (IOBs) provide predictive capabilities that traditional signature-based detection methods cannot match. Al-powered anomaly detection methods combine statistical methods with machine learning to create dynamic baselines that adapt to changing user and system behavior while reducing false positives through contextual analysis.

Zero Trust architectures customized for Al systems provide comprehensive protection throughout the Al lifecycle. Zero Trust Al Access (ZTAI) principles include strict authentication for all interactions with Al systems, least-privilege access to models and data, and continuous monitoring of Al behavior and data flow.

Quantum-resistant approaches prepare organizations for post-quantum threats through hybrid cryptography that combines classical and quantum-safe algorithms. Al-powered algorithm development for post-quantum cryptography uses machine learning optimization to improve both security and performance characteristics.



Metrics and ROI justify strategic investment

Comprehensive analyzes of purple teaming programs show measurable business value that supports executive decision making. Organizations that implement mature Purple Team programs achieve a 3:1 to 5:1 return on their security investment, with measurable improvements in threat detection rates, incident response times and overall security effectiveness.

ETRICS & RO

Purple Teams 427%

Purple team programs deliver 427% average ROI over three years SecOps 62%

Security operations efficiency improves by 62% through purple teaming Insurance

23%

Organizations with mature purple teaming capabilities command 23% higher cyber insurance coverage at 38% lower premiums

Performance indicators provide a clear framework for measurement. Threat Resilience Metrics show the percentage of techniques detected and blocked by MITRE ATT&CK, while Al-specific Detection Metrics measure the detection accuracy of adversarial samples and the prevention rate of model poisoning attacks. Operational metrics show the frequency of the purple team's exercises, the speed of the blue team's improvements and the effectiveness of cross-functional collaboration.

Cost-benefit analyzes help justify the budget through quantified risk mitigation. Industry benchmarks show that the cost of purple team programs ranges from \$200,000 per year for small organizations to more than \$3 million for large companies, with ROI typically negative in the first year due to set-up costs, but a 300-500% return on investment in the third year due to continuous improvement.

Financial justification models address executives' concerns about measuring the value of cybersecurity investments. Business impact frameworks quantify maintaining customer trust, protecting brand reputation and avoiding fines. Strategic benefits include improved competitive position, reduced friction in AI adoption and better talent retention through advanced skills development opportunities.



Implementation challenges reveal practical solutions

The analysis of successful implementations identifies common patterns and practical solutions to organizational challenges. Technical integration challenges focus on tool compatibility and data correlation, which are solved through the standardization of MITRE ATT&CK frameworks, automation platforms and integrated security operations centers.

RACTICAL

Purple Teaming 67%

67% of purple team initiatives fail within the first year due to organizational resistance

Purple Teaming 91%

Purple Teams with executive sponsorship achieve 91% success rates \$4.48m

The average cost of Al-enhanced phishing breaches now reaches \$4.88 million

The biggest obstacle to the success of purple teams is the management of organizational change. The competitive dynamics of red and blue teams require clear communication about collaboration goals, shared success metrics and regular cross-team training. Support from senior management proves essential to overcome cultural resistance and ensure sustainable resource allocation.

Staffing models adapt to organizational requirements and resource constraints. Dedicated purple team structures are suitable for large organizations with mature security operations, while virtual team models allow mid-sized organizations to achieve the benefits of collaboration without full-time resources. Hybrid approaches strike a balance between specialized expertise and broader organizational participation.

The most successful implementations start early in Al development lifecycles, invest in hybrid talent that combines offensive security with Al/ML knowledge, and establish continuous improvement processes based on measurable results. Organizations that view purple teaming as an ongoing capability rather than a periodic exercise achieve much better results.



A new methodology framework emerges

The synthesis of research findings reveals clear requirements for a next-generation AI security methodology that addresses the current limitations of the framework while leveraging proven purple teaming approaches. The proposed Strategic AI Defense Integration (SADI) framework provides a practical implementation guide that combines strategic risk management with operational security capabilities.

Adaptive Risk Assessment Engine that enables automated, continuous assessment of AI security risks based on organizational context, threat data and regulatory requirements. This closes the gaps in MITRE ATLAS' risk prioritization and provides managers with clear decision-making tools.

The Implementation Playbook Library provides industry-specific, step-by-step security implementation guides that translate strategic frameworks into operational capabilities. This links the governance approach of the NIST AI RMF to practical technical controls and measurable results.

The Executive Decision Dashboard provides a real-time visualization of the AI safety posture specifically designed for executive communication and strategic decision making. This overcomes the technical complexity barriers of current frameworks while supporting budget justification and program control.

Continuous Purple Teaming integration that embeds Al-specific joint safety testing throughout the Al lifecycle. This extends traditional Purple Teaming methods with automated tools, Alspecific attack simulations and continuous improvement processes.

Regulatory compliance automation that enables continuous monitoring and reporting for evolving AI safety regulations in multiple jurisdictions. This addresses the complexity of global compliance requirements while reducing organizational overhead and ensuring proactive adaptability.

Supply Chain Security Orchestration, which integrates AI workflow protection, vendor risk management and third-party security assessment into unified governance frameworks. This addresses the unique challenges of AI security in the supply chain while providing practical guidance for implementation.

The SADI framework leverages the strengths of existing frameworks while closing their implementation gaps through automation, executive engagement tools and practical guidance that organizations can implement immediately and continuously improve.





Organizations must act decisively to combat Al-powered threats while the regulatory window is still open. The convergence of rapidly proliferating Al attacks, evolving compliance requirements and proven purple teaming effectiveness creates both urgent risks and strategic opportunities.

Immediate actions for the next 90 days should include a comprehensive inventory of Al assets to understand current exposure, the establishment of Al governance structures with clear executive accountability, and the implementation of basic security controls, including API protection and access monitoring for existing Al systems.

Near-term initiatives for the next year will require investment in AI security platforms with automated threat detection capabilities, increased developer training for MLOps security practices, and the adoption of AI-specific incident response procedures that take into account the unique characteristics of AI system compromises.

Long-term strategic positioning through 2027 requires the deployment of advanced AI redteaming capabilities, preparation for evolving regulatory frameworks, including compliance with the EU AI Act, and the development of quantum-resistant security architectures that will become critical as quantum computing capabilities mature.

Success requires treating AI security as a strategic imperative rather than a technical concern. Organizations that embed security throughout their AI journey, from initial research and development through production deployment and ongoing operations, will achieve competitive advantages through enhanced trust, reduced risk exposure, and operational resilience.

The research demonstrates that Al-enabled threats represent fundamental shifts in cybersecurity risk that cannot be addressed through incremental improvements to existing approaches. Organizations that proactively implement comprehensive Al security capabilities through proven purple teaming methodologies, supported by executive leadership and measured through clear business metrics, will be best positioned to realize Al's transformational benefits while managing its inherent risks.

PRATEGY

Al Security \$4.7m

Companies that implement AI security measures within the next 12 months will save an estimated \$4.7 million

Valuation 34%

First-mover advantage in AI security translates to 34% higher market valuation **Purple Teaming**

7.2%

Every month of delay in implementing purple teaming increases breach probability by 7.2%

Contact Us



Contact Information

Sales sales@pentrasecurity.com