# Can You Trust AI in the Cloud?

## Exploring the Risks of Automated Security Analysis

A critical examination of AI-powered security tools in cloud-native environments and the trust boundaries we must establish.

**CornCon 2025 | Oct 10th, 2025**

# About Me

## Advait Patel

### Senior Site Reliability Engineer at Broadcom

Specializing in Zero Trust security, cloud infrastructure, and DevSecOps with over 8 years of experience securing cloud-native environments and AI workloads.

- **IEEE Leadership:** Region 4 Large Section Support Chair
- **Published Author:** IAM & Security for Google Cloud Platform
- **Industry Contributor:** OWASP GenAI Security founding member
- **Creator of DockSec:** AI-powered Docker Security Analyzer
- **Book Author:** Books on AI, Cloud, and Security published with Wiley, Springer Nature, and Apress
- **Conference Speaker:** Regular speaker at SANS Conference, Blue Team Con, OWASP Global AppSec etc on AI in security
- **Technical Articles & Publications:** Contributor to various industry journals and online platforms, including ISACA, ISSA, CSA, Hackernoon, The New Stack

# The AI Security Revolution

### Intelligent Scanning

AI-powered vulnerability detection analyzes cloud configurations at scale, identifying security gaps faster than traditional tools.

### Automated Remediation

GPT-driven recommendations suggest fixes for security issues, promising to reduce manual intervention and accelerate response times.

### Predictive Analysis

Machine learning models predict potential threats by analyzing patterns across cloud environments and threat intelligence feeds.

Organizations are racing to adopt these capabilities, driven by the promise of faster detection, reduced security debt, and scalable protection across complex cloud-native ecosystems. But speed without accuracy creates new vulnerabilities.

# Can we truly trust AI to secure our cloud environments?

This isn't a theoretical concern. As AI tools gain autonomy in security pipelines, the stakes rise dramatically. A misclassified threat, a hallucinated fix, or an overlooked vulnerability can have cascading consequences in production systems.
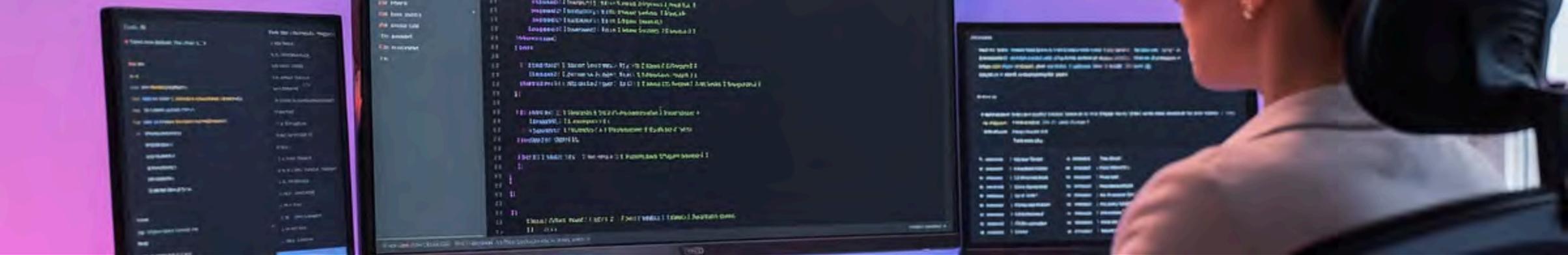
# The Trust Paradox

## Why Organizations Trust AI

- Processes massive data volumes instantly
- Operates 24/7 without fatigue
- Learns from vast threat databases
- Reduces time-to-detection dramatically
- Scales across distributed environments
- Provides consistent analysis

## Why We Should Be Skeptical

- Training data may contain biases
- Context understanding remains limited
- Hallucinations can suggest dangerous fixes
- Black box decision-making lacks transparency
- Edge cases often misclassified
- No accountability mechanism built-in

The tension between AI's capabilities and its limitations creates a critical gap in modern cloud security strategies. Understanding this paradox is the first step toward responsible implementation.

# Introducing DockSec

To understand these risks concretely, let's examine **DockSec**—an open-source, AI-powered Dockerfile security analyzer. This tool represents the current generation of automated security analysis: fast, intelligent, and deeply integrated into DevSecOps pipelines.

### What It Does

DockSec analyzes Dockerfiles for security vulnerabilities, misconfigurations, and best practice violations. It uses machine learning to understand context and provide actionable remediation guidance.

### The Promise

Automated security checks at commit time, catching issues before they reach production. Developers receive immediate feedback without waiting for security team review.

### The Reality

Like all AI tools, DockSec makes mistakes. Some are minor. Others could be exploited in production environments, creating a false sense of security.

# Case Study: When AI Gets It Wrong

## A Real-World Misclassification

Let's walk through a specific example where DockSec incorrectly identified a Dockerfile security issue. This isn't just an academic exercise—this pattern of failure appears across AI security tools and has real consequences.

### 1 The Scenario

A development team uses a legitimate base image from a trusted registry. The Dockerfile follows standard practices for layer optimization and includes proper USER directives.

### 2 AI Analysis

DockSec flags the configuration as "high risk," claiming the base image is outdated and recommending an immediate switch to a different image—one that's actually less maintained.

### 3 The Impact

Following the AI's recommendation would introduce actual vulnerabilities. The suggested image has known CVEs and lacks the security hardening present in the original choice.

# Anatomy of the Failure

### 1

### Context Blindness

The AI analyzed image age without understanding the organization's update schedule or compatibility requirements. It lacked context about why this specific version was chosen.

### 2

### Hallucinated Alternatives

The recommended "secure" alternative didn't exist in the specified registry. The AI generated a plausible-sounding but incorrect image name based on pattern matching.

### 3

### Overconfident Scoring

The tool assigned a 95% confidence score to its incorrect analysis, making the false positive appear highly reliable to users who trust the metrics.

**Key Insight:** AI confidence scores don't reflect actual accuracy. High confidence can accompany completely incorrect analysis, especially when the model encounters edge cases outside its training data.

# How This Gets Exploited

Security misclassifications aren't just operational problems—they create attack vectors. Adversaries actively exploit the trust organizations place in AI security tools.

## Attack Pattern #1: False Negatives

Attackers craft payloads specifically designed to evade AI detection. By understanding model behavior, they introduce malicious configurations that the AI classifies as safe.

- Obfuscated command patterns
- Legitimate-looking but dangerous configurations
- Exploiting known model blind spots

## Attack Pattern #2: Poisoned Recommendations

More sophisticated: manipulating the AI to recommend vulnerable configurations. This turns the security tool into an attack vector itself.

- Suggesting outdated dependencies
- Recommending insecure alternatives
- Creating backdoor opportunities

# The Blind Spots in DevSecOps Pipelines

## Development

AI scans code at commit time, flags false positives, trains developers to ignore warnings

## Build

Automated builds trust AI "pass" results without human verification of critical paths

## Testing

Security tests assume AI pre-screening caught major issues, reducing coverage

## Monitoring

Runtime monitoring relies on AI threat detection with similar blind spots

## Deployment

Production systems inherit undetected vulnerabilities from earlier stages

The integration of AI throughout the pipeline creates compounding risk. Each stage assumes the previous stage's AI analysis was accurate, but errors propagate silently through the entire system.

# When to Trust AI

AI isn't inherently untrustworthy—it's a tool that excels in specific contexts and fails in others. The key is understanding these boundaries and implementing appropriate oversight.

### Pattern Recognition at Scale

AI excels at detecting known vulnerability patterns across thousands of files. Trust it for initial screening of common misconfigurations and well-documented CVEs.

### High-Volume Triage

Use AI to prioritize security alerts, ranking them by likely severity. This reduces noise and helps security teams focus on genuine threats requiring investigation.

### Consistency Checking

AI reliably catches policy violations and configuration drift when rules are clearly defined. It won't forget to check required fields or skip validation steps.

# When to Intervene

### Critical Infrastructure Decisions

Never allow AI to make autonomous decisions about production infrastructure, authentication systems, or network security boundaries. These require human judgment and accountability.

### Context-Dependent Analysis

When security decisions depend on business logic, compliance requirements, or organizational context, AI lacks the necessary understanding. Human expertise is essential.
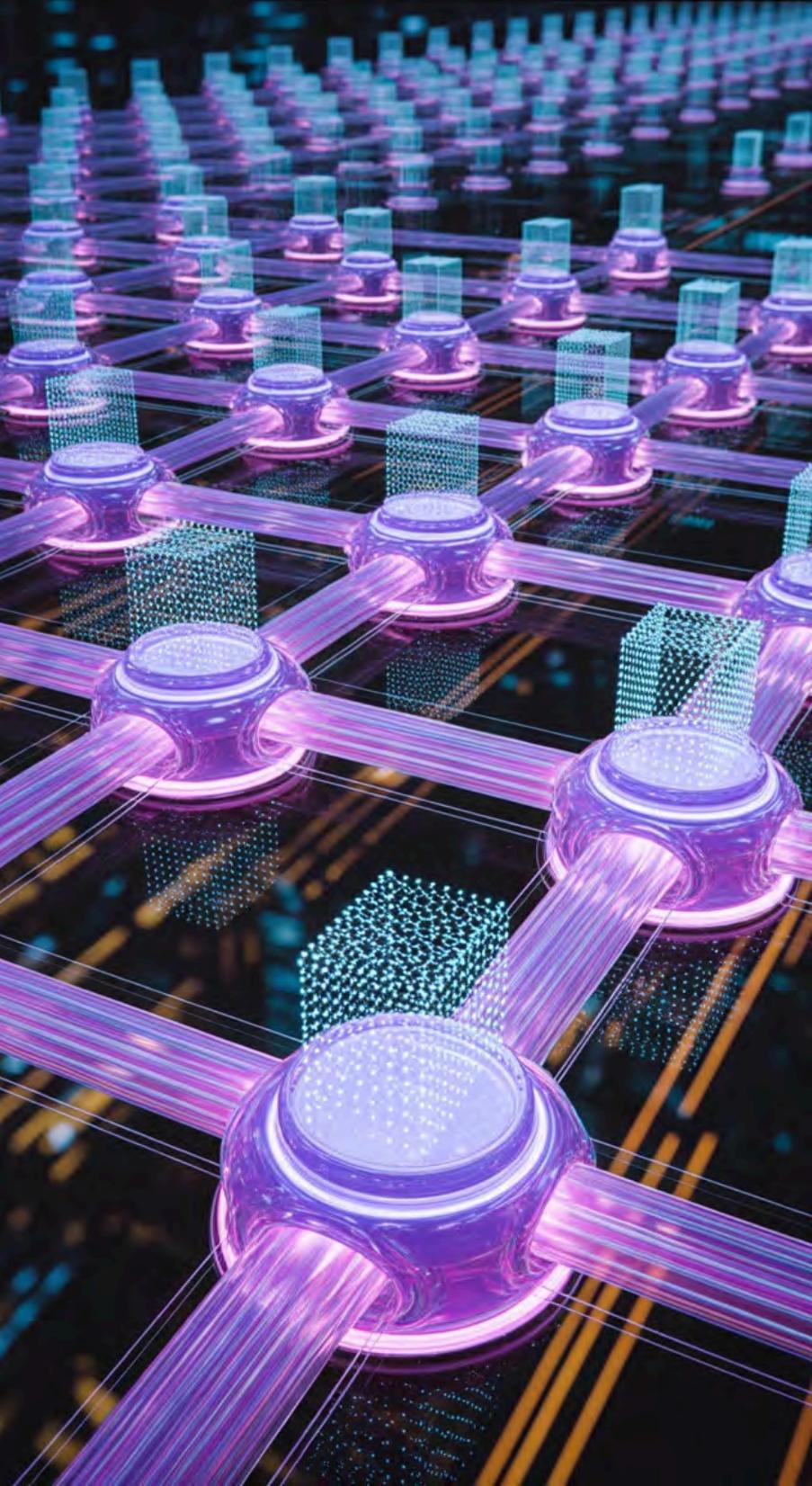
### Novel or Unusual Configurations

AI performs poorly on edge cases outside its training data. Custom solutions, experimental technologies, or unique architectural patterns require manual security review.

### High-Confidence False Positives

When AI shows high confidence in findings that seem questionable, always investigate. Overconfident misclassifications are a red flag for model limitations.

# Building Effective AI Guardrails

Responsible AI implementation requires deliberate safeguards that combine automation with human oversight. These guardrails prevent AI failures from becoming security incidents.

## Layer 1: Input Validation

Verify that AI tools receive clean, representative data. Implement sanity checks on recommendations before they reach decision points.

## Layer 2: Confidence Thresholds

Require human review for findings below 80% confidence or above 98% confidence (both indicate potential issues). The middle range is most reliable.

## Layer 3: Change Control

All AI-recommended security changes must go through standard change management. No autonomous remediation in production environments.

## Layer 4: Continuous Validation

Regularly audit AI decisions against ground truth. Track false positive and false negative rates. Retrain or replace models that drift.

# A Practical Governance Framework

## Risk Assessment

- Classify AI tools by criticality
- Map potential failure modes
- Define acceptable error rates
- Establish escalation paths
- Document known limitations

## Implementation Controls

- Mandatory human review gates
- Dual-control for critical changes
- Audit logging of all AI decisions
- Version control for AI models
- Regular penetration testing

## Ongoing Governance

- Quarterly accuracy reviews
- Incident response for AI failures
- Training on AI limitations
- Vendor security assessments
- Compliance documentation

**Critical Success Factor:** Governance frameworks must be living documents. AI capabilities and risks evolve rapidly—your oversight mechanisms must evolve with them.

# The Future of Cloud Security

AI will continue transforming cloud security, but success requires a fundamental shift in how we think about these tools. They're not replacements for human expertise—they're force multipliers that work best under informed human guidance.

### Trust Through Transparency

Demand explainable AI from vendors. Security decisions require clear reasoning, not black-box outputs. Build internal expertise to evaluate AI recommendations critically.

### Defense in Depth

Never rely solely on AI security tools. Layer multiple detection methods, maintain manual review processes, and assume AI will occasionally fail catastrophically.

### Continuous Learning

Both humans and AI must learn from failures. Create feedback loops, share lessons across teams, and maintain a culture of healthy skepticism toward automation.

# Thank You & Q&A

Questions? Let's discuss how we can build more trustworthy AI security systems together. This conversation doesn't end here—it's just beginning.

**Contact:**
advaitpa93@gmail.com
LinkedIn: /in/advaitpatel

**Session Resources:**
CornCon 2025

## Your Questions, Our Insights.