

# AVAST!

## Your AI Is Writing Vulnerable Code

A Threat Modeling Framework for AI-Assisted Development

Brad Tenenholtz | CornCon 2025 | October 10-11, 2025

RiverCenter, Davenport, Iowa

# Agenda

## **The Current State**

AI code vulnerability statistics and adoption gaps

## **The Problem Space**

Where vulnerabilities happen and why traditional tools fail

## **Threat Modeling Fundamentals**

Why threat modeling is essential for AI code security

## **AVAST Framework**

A specialized threat model for AI-assisted development

## **Implementation**

Organizational transformation and Monday morning actions

# The Data Bomb & Adoption Gap

## Security Check Failure Rate



- 45% of AI-generated code fails basic security tests
- 62% contain formally verifiable vulnerabilities
- 72% failure rate for Java

## Adoption & Concerns

**83%**

of organisations use AI to generate code

**92%**

of security leaders express concern

**63%**

lack visibility & consider banning AI

# Better Models are Not More Secure

## Security Pass Rate

Large Models (>100B)

51%

Medium Models (20-100B)

51.5%

Small Models (<20B)

50.5%

GPT-4o (1.7T) (Baseline)

51%

## Security doesn't scale with model size

Huge models like GPT-4o (1.7T params) produce code as vulnerable as models 100x smaller.

# Where AI Vulnerability Happens

## OWASP Bugs

AI rarely implements proper input sanitization, leading to injection vulnerabilities. Paradoxically, SQL injection vulnerability by AI code has actually fallen in the past year.

## Secrets Management

Agentic rules files (.cursorrules, .claude, etc) often contain secrets. AI-based CLI tools that scan secrets may regurgitate them in other codebases. These "amplifier" attacks are a novel threat.

## Authorization and logging

Mishandled sessions, poor or inappropriate logging, and other architectural problems are an AI hallmark

```
// Common AI-generated anti-pattern
try {
  // Complex operation
  const result = await database.query(`SELECT * FROM users WHERE id = ${userId}`);
} catch (error) {
  console.log("Error:", error); // Exposes stack trace
  return null; // Fails open
}
```

**Key Finding:** AI consistently generates the same vulnerable patterns across different models and contexts

# The Iteration Paradox

More AI iterations = More vulnerabilities

## Why This Happens:

- ▶ Each iteration adds complexity without security context
- ▶ AI optimizes for functionality, not security
- ▶ Technical debt compounds with each generation
- ▶ Security assumptions get lost in translation

## Real Example:

**Iteration 1:** Basic authentication

**Iteration 2:** Added "remember me" (stored plaintext)

**Iteration 3:** Added social login (no CSRF protection)

**Iteration 4:** Added API keys (hardcoded in client)

Each "improvement" added new vulnerabilities

After three unsuccessful AI prompts, clear history and start over.

# The Next Threat: Poisoned AI Sources

## EchoLeak Attack Pattern

Malicious actors are poisoning training data to make AI generate vulnerable code

### Current Threats:

- ▶ Deliberately vulnerable GitHub repos with high stars
- ▶ Poisoned Stack Overflow answers
- ▶ Malicious MCP (Model Context Protocol) servers
- ▶ Compromised documentation sites

### Attack Vector:

```
// Looks legitimate, subtly vulnerable
function authenticate(user, pass) {
  // "Optimized" query
  return db.query(
    `SELECT * FROM users WHERE
    username='${user}' AND
    password='${pass}'`
  );
}
```

### The Amplification Effect:

One poisoned example → Thousands of AI suggestions → Millions of vulnerable applications

# Why Our Security Tools Often Aren't Enough

## Pattern Matching

SAST, DAST, Port scanning, all produce noncontextual findings. With AI, these findings will be overwhelming to deal with.

*Without context, scanners are useless.*

## Architecture Review

Reviewers assume code was written with intent and context

*AI code lacks security reasoning*

## Security Training

Teaches developers, not AI. Teaching AI is possible, but requires a skillset we are not yet developing.

*AI Centers of Excellence are critical to control and review agentic rules files.*

## The Fundamental Problem:

We have known for years now that we have to bring context to our security work in order to be effective.  
What AI requires of us is nothing new - it's a forcing function to drive us to the modernization we already needed.

## What We Need:

- ▶ Pattern detection for AI-generated code
- ▶ Context validation for security assumptions
- ▶ **Threat modeling specifically for AI development**

# What is Threat Modeling?

"Threat modeling is a structured approach to identifying and prioritizing potential threats to a system"

## Why We Need It for AI Development:

### Traditional Development

- ▶ Developer understands context
- ▶ Security is considered during design
- ▶ Threats are anticipated
- ▶ Mitigations are intentional

### AI-Assisted Development

- ▶ AI lacks context awareness
- ▶ Security isn't part of generation
- ▶ Threats are copied from training data
- ▶ Mitigations are often missing

## The Value of Threat Modeling:

- ▶ Identifies vulnerabilities before they're coded
- ▶ Focuses security efforts where they matter most
- ▶ Provides context AI lacks
- ▶ Creates testable security requirements

## Example: Authentication System Threat Model

Threat	Vulnerability	Mitigation	Action Required?
SQL Injection in login	(SAST) Not all queries are properly sanitized	Ensure all queries properly use stored procedures	Yes - Parameterized queries required
Attacker compromises gateway server to redirect requests	(Port scanner) Known RCE vuln on OS	Update OS	No – gateway is not accessible to the Internet
3rd party library poisoning	(SCA) 3rd party libraries in use	Whitelisted libraries	Yes – secure 3rd party libraries
Takeover of SQL Database server	(CNAP) Misconfigured cloud resources	Secure IaC	No – only stores pictures of cats

### Key: Tools produce findings, not vulnerabilities:

**A finding is a vulnerability only when we can match it to a threat:** If it isn't vulnerable, don't fix it!

Result: Spend more time on fewer, higher impact issues.

# Applying the Threat Model

## Before AI Generation:

```
/* SECURITY REQUIREMENTS FOR AI:  
 * 1. Use parameterized queries only  
 * 2. Hash passwords with bcrypt (min 10 rounds)  
 * 3. Include rate limiting (5 attempts/minute)  
 * 4. Regenerate session on login  
 */  
  
// Generate login function here...
```

## After AI Generation:

### ✓ Check Against Model

- ▶ Parameterized queries used?
- ▶ Bcrypt implemented?
- ▶ Rate limiting present?
- ▶ Session regeneration?

### ⚠ Common AI Misses

- ▶ Used string concatenation
- ▶ Improper logging
- ▶ Strange resource utilization and fallbacks
- ▶ Session not regenerated

# Introducing AVAST

A More Precise STRIDE for AI Coding Use-Cases

## Traditional STRIDE

Spoofing  
Tampering  
Repudiation  
Information Disclosure  
Denial of Service  
Elevation of Privilege

Generic threats for any system

## AVAST for AI Code

Authentication Flaws  
Validation Gaps  
Auditing  
Secrets in Code  
Trust Boundaries

Specific to AI-generated vulnerabilities

**AVAST: Designed for the next generation of vulnerabilities**

# The Organizational Transformation

To implement threat model-centric AI security:

## 1. Awareness

Train developers on AI-specific vulnerabilities

Show them what AI gets wrong consistently

## 2. Integration

Add threat modeling to AI development workflow

Before generation, not after deployment

## 3. Tooling

Deploy AI-aware security scanning

Detect patterns, not just vulnerabilities

## 4. Governance

Establish AI code review requirements

Different bar for AI vs human code

## The Cultural Shift:

**From:** "AI makes us faster"

**To:** "AI makes us faster when properly governed"

**Key Success Factor:** Make security easier than insecurity

Provide templates, tools, and training that make secure AI development the path of least resistance

# Monday Morning Actions

## Step 1: Discover

```
git grep -l  
"copilot\|cursor\|claude" --all
```

Find your AI coding footprint

## Step 2: Assess

Run AVAST checklist on recent AI code:

- ▶ Authentication flaws?
- ▶ Validation gaps?
- ▶ Auditing and logging?
- ▶ Secrets in code?
- ▶ Trust assumptions?

## Step 3: Implement

Start with one team:

- ▶ 15-min AVAST training
- ▶ Template for AI prompts
- ▶ Pre-commit security check
- ▶ Weekly metrics review

## The 30-Day Goal:

Reduce AI-generated vulnerabilities by 50%

Measurable. Achievable. Essential.

# Works Cited

**Veracode.** (2025). *GenAI Code Security Report: Assessing the Security of Using LLMs for Coding*. Retrieved from [veracode.com/genai-report-2025](https://veracode.com/genai-report-2025)

**GitHub.** (2024). *The State of AI-Assisted Development*. GitHub Octoverse Report.

**MITRE.** (2024). *Common Weakness Enumeration: AI-Generated Code Patterns*. CWE-2024-AI.

**Stanford Security Lab.** (2024). *EchoLeak: Poisoning LLM Training Data for Targeted Vulnerabilities*. arXiv:2024.13579.

**Microsoft.** (2024). *STRIDE Threat Modeling*. Microsoft Security Development Lifecycle.

## Questions?

Brad Tenenholtz

[brad.tenenholtz@gmail.com](mailto:brad.tenenholtz@gmail.com)

<https://www.linkedin.com/in/the-real-brad-tenenholtz/>