# Use of AI is Here to Stay: Enabling Innovation Responsibly and Securely



**Todd Hathaway**

*Global Head of AI Security
and Cyber Innovation*

www.linkedin.com/in/toddhathaway



World Wide Technology

World Wide Technology

BLOG

RSA and the Agentic AI Bandwagon

RSA Conference 2025
RESPONSIBLE SECURE AGENTIC AI
RSA CONFERENCE 2025

RSAC 2025 Conference
San Francisco · April 28 – May 1 · Moscone Center
Many Voices. One Community.
RSAC™ 2025 Keynote
Security in the Age of Agentic AI
Vasu Jakkal
Corporate Vice President
Microsoft Security

GenAI SECURITY PROJECT
https://genai.owasp.org
Join Us Live!
Agentic Security Open Workshop Livestream Event from RSAC 2025
April 30th, 2:00–5:00pm PST

metron security | RSAC 2025 Conference
Trending Now at RSAC:
Agentic AI

Autonomous TPRM with Agentic AI
Find out more at RSAC 2025 Conference
San Francisco · April 28 – May 1 · Moscone Center

RSAC 2025 Conference
The First Agentic AI Security Champion for ASPM
ArmorCode

RSAC
2025 Conference
Agentic AI, Global Recognition, and the Cowboys of GenAI Security

RSA Conference 2025
The dawn of agentic AI in security operations

A2A  MCP

World Wide Technology

# AI is prevalent everywhere

**75%** [1]
Employees Use AI

**69%** [2]
Developing AI Applications

**50%** [3]
Will Use AI Agents by EOY

% of Enterprises

85%
80%

75%

69%

50%

10%

2023    2024    2025    TODAY    EOY 2025

[1] EY, EY 2024 Work Reimagined Survey
[2] Wiz, State of AI in the Cloud 2024
[3] Capgemini, Generative AI in organizations 2024

World Wide Technology

# AI is prevalent everywhere – Risk vs. Value

World Wide Technology

**Risk**

**Value**

Agentic AI

Enterprise Developed AI

Usage of AI Apps

- Business Disruption
- Data Breach

- Business Logic Manipulation
- Model Theft
- Data Poisoning

- Overreliance
- Oversharing
- Data Loss

- Knowledge Retrieval
- Writing & Content Generation
- Code Generation

- Processing Unstructured Data
- Business Intelligence
- Fraud Prevention

- Business Automation
- Cost Reduction
- Force Multipliers for speed & efficiency

# Security means different things

**Securing AI Systems**

**Using AI Securely**

**Defending from Adversarial AI**

**Using AI for Cyber**

# Domains of AI Security



## Adversarial Use of AI

- Deepfakes & Misinformation
- Phishing
- Malware
- Social Engineering
- Denial of Service
- Surveillance & Espionage
- AI Poisoning
- API Reconnaissance
- Credential Stuffing

## Governance, Risk and Compliance

| | |
|---|---|
| Program Strategy Development | Program Security Maturity Assessments |
| Policies and Procedures | Awareness Training |
| Controls Gap Assessment | Model Risk Management |
| Compliance Measurement | Data Governance & Classification |

## Security of AI Systems

- Data Readiness, Security & Privacy
- Model Scanning & Model Theft
- Guardrails/Firewalls
- API & Agentic Services
- Red Teaming & AI-SPM

## Security of AI Usage

- Usage Discovery
- Data Loss Protection
- 3rd Party Model Risk Management
- Agentic Tools & MCP,A2A,AGNTCY
- Regulatory

## AI for Cyber Security

- Risk Quantification & Compliance
- Threat/Vulnerability Management
- Security Ops (SIEM/SOAR)
- Identity & Access Management
- AI Code Remediation
- Fraud Detection
- App Detection & Response
- Endpoint Detection & Response
- Threat Detection & Response

Trustworthy    Responsible    Ethical    Accountable    Transparent    Secure

World Wide Technology

# Current Trends in AI Security

Agentic….MCP…A2A

**World Wide Technology**

## Organizational Areas of Focus for AI Security

### Security of AI Systems

- Data Readiness, Security & Privacy
- Model Scanning & Model Theft
- Guardrails/Firewalls
- API & Agentic Services
- Red Teaming & AI-SPM

### Security of AI Usage

- Usage Discovery
- Data Loss Protection
- 3rd Party Model Risk Management
- Agentic Tools & MCP
- Regulatory

### Adversarial Use of AI

- Deepfakes & Misinformation
- Phishing/Social Engineering
- Malware
- API Reconnaissance
- BOT/ATO/DOS

## Current Industry Trends

- Recognition as a standalone topic
- Agentic MCP vs. A2A
- Authorization for data sources
- Data security for AI Agent access
- Red Teaming of AI Systems
- AI GW vs. FW vs. Guardrails
- Shadow AI --→ AI Everywhere
- Secure by Design for AI Apps & Agents

## Adversarial Use of AI

- Phishing
- Vulnerability Exploitation
- CAPTCHA Breaking AI
- Vulnerable API Reconnaissance
- AI-Generated Content including deepfakes and malicious code

## Vendor Landscape

- Startup activity is Still accelerating
- Acquisition activity increasing
- Lack of standard terminology creating confusion

## Barriers to Adoption

- Governance Delays
- Privacy Concerns
- Explainability/Auditability

Discussions of artificial intelligence (AI) often swirl with mysticism regarding how an AI system functions.

The reality is far simpler:

"AI is a type of software system."

– CISA

World Wide Technology

**Newsweek**

AI | Hospital   Health Care   Artificial Intelligence   CEO   Security   Cybersecurity   Data   Leak   Safety

# Average Health System Audit Finds 70 'Quiet' AI Applications, CEO Says

Published Apr 07, 2025 at 2:51 PM EDT

# Security of AI Usage

What AI tools are in use in your enterprise?

**World Wide Technology**

# AI impacts every corporate persona

**Every role & department has a need for AI tools**

| Legal | Finance | Sales | R&D/IT | Marketing | HR |

*78% of AI users are bringing their own AI tools to work*

– 2024 Work Trend Index Annual Report from Microsoft

Everyone is using AI.  AI is being applied to every industry.
YES, you can embrace the innovation while you manage the risk

# AI comes in many forms

Websites

Mobile Apps

Embedded AI (SaaS)

Chatbots

AI Assistants

Copilots

Voice Assistants

AI Coding Assistants

Applications

AI-enabled Hardware

AGENTIC AI

Agentic AI can Autonomously:
**Perceive - Reason - Act**
To achieve desired outcomes with minimal human intervention.

# Example of Enterprise AI Tool Landscape

# GenAI Usage Data from real customers (Q1, 2025)

**World Wide Technology**

**176k**
Prompts

**8k**
Users

**8.2k**
Files

**254** — Average Number of Apps in Use in each enterprise

**45.4%** — Of sensitive data submissions were using accounts via personal accounts

**7%** — Of users accessed Chinese-based apps with data training enabled (DeepSeek, Manus, Ernie Bot, Qwen Chat, Baidu Chat)

# Enabling AI Usage Securely

Discovery & Inventory → Policy & Guidance → Assess & Approve → Enforce & Protect

Discovery & Inventory → Policy & Guidance → Enforce & Protect → Coach & Retrain

# Control Elements for AI Usage Security

World Wide Technology

| AI Usage Assessment | **Continuous Discovery & Inventory** |
|---|---|

| Acceptable Use Policy | **Continuous Risk Assessment** |
|---|---|

| User Education | **Enforce & Protect** | User Education |
|---|---|---|

| Allow/Block? | Data Protection | Advise & Coach |
|---|---|---|

| File Upload | Redact or Obfuscate |
|---|---|

| GPTs | LLM Apps | Ai in SaaS | AI Agents | Copilots | AI Code Gen | Mobile Apps |
|---|---|---|---|---|---|---|

# Tool Approaches to AI Security



**Solution Approaches:**

**Less Complex**

✓ API Connectivity

✓ Browser Extension / Plug-in

✓ Agent / SSE Solutions

✓ Dedicated Browser

✓ Proxy / Gateway / Firewall

✓ Discovery

✓ MDM

✓ Combination of the above

**More Complex**

- API / Logs
- Cloud Agent

SSE

Public Cloud     SSE     Private Cloud

Email, Records and Compliance Documentation

Corporate Router     Network Printer

Discovery     Proxy / Gateway / Firewall

Administration Team

- MDM
- Dedicated Browser

- Browser Extension / Plug-in
- Dedicated Browser
- Agent

# Why network visibility is not enough



**Reverse Proxies -** Cannot prevent data exposure on unsanctioned apps

**MITM Decryption –** Scaling Limitations, Performance Limits, PQE Impacts looming?

**API Scanning** Cannot prevent malicious activity within sanctioned apps

**Forward Proxies** Cannot provide access control on unmanaged devices

# AI Agents vs. GenAI vs Agentic AI

World Wide Technology

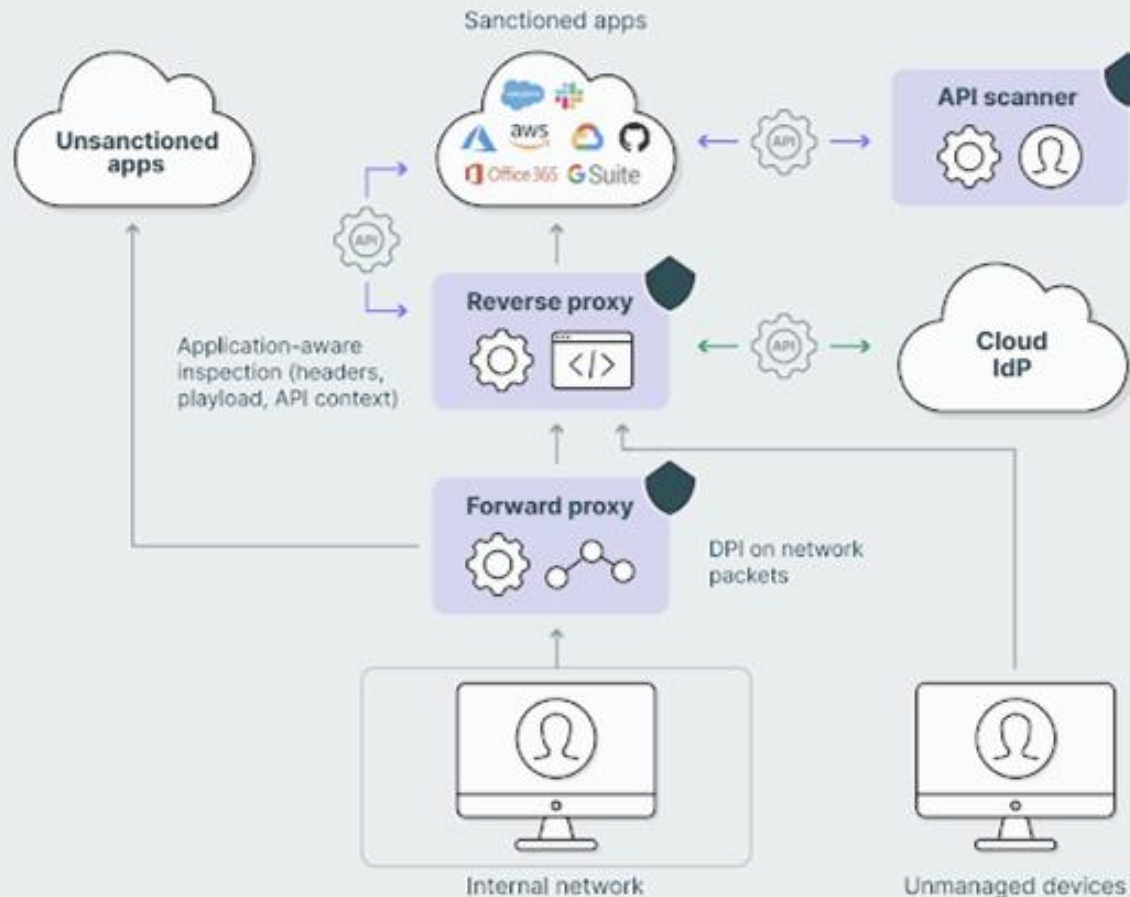|  | AI Agents | Gen AI | Agentic AI |
|---|---|---|---|
| **Functionality** | Automated task execution based on rules or patterns | Content Creation based on training, patterns, & predictions | Autonomous Action, Problem Solving and Decision-Making |
| **Adaptability** | LOW-Follows fixed workflows | MED-Can generate varied responses | HIGH-reasons, adapts, plans, and acts |
| **Business Use** | Automating repetitive tasks | Read and Summarize, or Generate text, images, or code | Optimize operational processes & make strategic decisions |
| **Examples** | Customer Service Bots, IT Automation | Marketing Content, Code Gen, Legal Assistants, Copilots, etc. | Autonomous Supply Chain, Cyber Analysts, Voice Agents |
| **Limitations** | Complex Reasoning | Needs accurate prompting and lacks independent action | Explainability and complete accuracy – Still need H-I-T-L |

A2A and the <u>Model Context Protocol (MCP)</u> are complementary standards for building robust agentic applications:

- **MCP (Model Context Protocol):** Connects agents to **tools, APIs, and resources** with structured inputs/outputs. Think of it as the way agents access their capabilities.

- **A2A (Agent2Agent Protocol):** Facilitates **dynamic, multimodal communication between different agents** as peers. It's how agents collaborate, delegate, and manage shared tasks.

# MCP

## Model Context Protocol

Defines a standardised interface for supplying structured, real-time context to **large language models.**

### CORE FUNCTIONALITIES

MCP lets you pull in external resources like files, database rows, or API responses - right into the prompt or working memory.

Rather than stuffing your prompt with every possible detail, MCP helps assemble just the context that matters.

MCP also lets models call tools dynamically.

---

## Application Host Process

**HOST**

**Client 1** — **Client 2** — **Client 3**

### Local Machine

**MCP Server 1**    **MCP Server 2**

**Local Resource**    **Local Resource**

### Internet

**MCP Server 3**

**Remote Resource**

**A2A**

Agent2Agent Protocol

Enable structured communication & coordination between **AI agents** operating in the environment.

**CORE FUNCTIONALITIES**

Facilitates message passing and task delegation between agents to coordinate actions.

Enables agents to share observations, goals, or partial outputs for collective decision-making.

Supports synchronization of agent states across distributed environments.

Agent

Local Agents

Vertex AI

Agent Development Kit

A2A

Agent

Local Agents

Vertex AI

Agent Development Kit

MCP

API's & Enterprise Applications

MCP

API's & Enterprise Applications

# A2A and MCP (more complementary than competitive)



https://google.github.io/A2A/#why-a2a-matters

https://google.github.io/A2A/topics/a2a-and-mcp/

# Early discoveries in MCP vulnerabilities



**Command Injection**

Using prompts with embedded meaning to trigger unauthorized MCP actions by the agents.

Effect: Moderate

**Server-Sent Events Problem**

SSE workflow creates latency and security issue due to constant opening of line during transfers

Effect: Moderate

**Server Data Takeover**

A compromised tool server can take over other servers data and passwords.

Effect: Severe

**Tool Poisoning**

Embedding malicious tools codes in MCP to manipulate their actions for a given task

Effect: Severe

**Privilege Escalation**

Malicious tools can override or intercept calls made to a trusted tools that you use.

Effect: Severe

**Persistent context**

MCP records your context through out your sessions, this can lead to context tampering

Effect: Small

# Wrap-Up

World Wide Technology

# Security of AI: Defining Product Categories

## AI Security Governance, Risk & Compliance

Software used to manage/govern enterprise AI assets and usage, assess & manage risk, and/or map controls to compliance requirements

## AI Discovery & Inventory

Software tools that provide automated discovery and inventory tracking for AI artifacts, 3rd Party AI Tools and 3rd Party AI embedded in existing SaaS. Many existing tools provide some capabilities here & can be used to begin and identify the need to add better controls)

## Security of AI Systems

| 1. Secure the Data | 2. Secure the Model | 3. Secure the Usage of the App |
|---|---|---|

### 1. Secure the Data

**Readiness & Risk Assessment for AI**
Is Data security ready for AI?

**Data Discovery, Classification, and Labeling** – Discovery of all the data that might be used in AI systems:, including structured and unstructured data source, Classification of all data types (e.g., PII, financial, IP), and Labeling of all data for appropriate use or protection

**Data Access Governance**
Ensures **only the right people, processes, and AI models** can access certain data and that models can only provide inferences to authorized data.

**Privacy Enhancing Techniques (PETs)**
Even if the AI needs sensitive data, privacy can still be protected.

Options include:
- **Data masking and tokenization**
- **Stochastic Randomization**
- **Homomorphic Encryption**

### 2. Secure the Model

**Supply Chain (Model Scanning, AI-BOM AI-SPM)**

This area is analogous with traditional Software Supply Chain security and covers the risks to code and artifacts in the AI/ML development lifecycle, including malicious code analysis, AI Bill-of-Materials, and AI Security Posture Monitoring, Source Dependencies, and provenance of artifacts from source to prompt .

**Red Teaming/Vulnerability Scanning**

Systematically probing both:
- **Models** that serve as central components for the applications, and
- **Systems** and **Data** used throughout the lifecycle of the application:

From model development and training, through application staging pipelines, and continuously in production runtime environments. Combines traditional adversarial testing with AI-specific methodologies, addressing risks like: **Prompt injection, Toxic outputs, Model extraction, Bias, Knowledge risks and Hallucinations.**

### 3. Secure the Usage of the App

**AI Runtime Security**

**Protecting AI systems *while they are running*.**
AI applications can be **tricked into leaking sensitive info, generate harmful or biased content**, and/or t**ake risky actions** based on bad inputs
**AI Firewall-**filtering inputs and outputs to block:
- Malicious prompts (like jailbreaks)
- Inappropriate or dangerous responses
- Unauthorized data access

**AI Guardrails** -rules and limits that keep the AI on track. They guide the AI to:
- Stick to **approved topics**
- Avoid **risky actions**
- Always follow **company policies**

Combined, the **AI Firewall** and **Guardrails** help enterprises ensure their AI behaves **safely, ethically, and within bounds** AI Runtime tools should have ability to filter both input and responses.

## Security of AI Usage

Security of Usage refers to any 3rd party AI usage and covers a broad area of security control types that can play a partial role in securing AI usage. Product areas include that offer security features for AI usage include:

**SASE/SSE**

**AI Focused (Browser Inserted)**

**AI Focused ( Network Inserted)**

**Enterprise Browser/Browser Extensions**

**Next Gen DLP**

**Digital Workforce Security (Security of AI Agents)**

**Code Development Focused**

**SSPM (SaaS Security Posture Management)**

Each Category has strengths and weaknesses and a complete program for using AI securely should start with a discovery exercise to gain a clear understanding of the current AI tools that are being used across the organization.

# Security of AI: WWT Market Landscape (v05122025)

World Wide Technology

## AI Security Governance, Risk & Compliance

CISCO · PROTECT AI · truyo · credo | ai · Tumeryk AI Trust Score · Pillar · CRANIUM · Holistic AI · fairnow · Enkrypt AI · preamble · FAIRLY · CitrusX

## AI Discovery & Inventory

### 1st Party Apps
CISCO · PROTECT AI · WIZ · truyo · Pillar · AppSOC AI

### 3rd Party Apps
CISCO · zscaler · paloalto NETWORKS · Acuvity · Aurascape · Singulr AI · WITNESS AI · onyx · unbound security · portal26 · SurePath AI

### 3P-SaaS
nudge · Reco · suridata

## Security of AI Usage

### SASE/SSE
paloalto NETWORKS · CISCO · zscaler · FORTINET · netskope · CHECK POINT

### GenAI Focused – (Browser Insertion)
Prompt · onyx · harmonic · AIM · Singulr AI · LASSO · Acuvity · Aurascape · magicmirror

### GenAI Focused – (Network MITM)
WITNESS AI · portal26 · Singulr AI · SurePath AI

### Enterprise Browser/Browser Extension
LayerX · SquareX · Island · paloalto NETWORKS · keep aware · SURF · MENLO SECURITY · primary · SERAPHIC

### Next-Gen AI-Focused DLP
cyberhaven · CYERA · harmonic · Nightfall AI · magicmirror · next

### Digital Workforce (Security of Agents)
onyx · Singulr AI · Acuvity · zenity

### Code Development Assistants
Prompt · Lumeus · LASSO

## Security of AI Systems

### 1. Secure the Data → 2. Secure the Model → 3. Secure the Usage of the App

#### 1. Secure the Data

**Risk Assessment**
KNOSTIC · opsin · Aurascape

**Discovery, Classification, Labeling**
BigID · CYERA · securiti · VARONIS · CROWDSTRIKE · paloalto NETWORKS · Microsoft · LightBeam.ai · RELYANCE AI · CONCENTRIC AI

**Data Access**
Polymer · KNOSTIC · Velotix · auth0 by Okta · securiti · privacera

**Privacy Enhancing Techniques**
PROTOPIA · privacera · ANTIMATTER · MIRROR Security · skyflow · PROTECTO · PRIVATE AI · inpher

#### 2. Secure the Model

**Supply Chain (Model Scanning, AI-BOM AI-SPM)**
PROTECT AI · paloalto NETWORKS · WIZ · Pillar · CISCO · REVERSINGLABS · HIDDENLAYER · NOMA

**Red Teaming/Vulnerability Scanning**
NVIDIA · splx AI · Lakera · PROTECT AI · MINDGARD · ADVERSA · CISCO · REPELLO AI · promptfoo · Calypso AI · NOMA · LASSO · GRAY SWAN · Javelin · Galileo · Virtue AI · TROJ.AI · dreadnode · TESTSAVANT · MIRROR Security · Enkrypt AI · HIDDENLAYER · Haize Labs

#### 3. Secure the Usage of the App

**AI Runtime Security**
NVIDIA · PROTECT AI · paloalto NETWORKS · CISCO · Akamai · f5 · CLOUDFLARE · Prompt · Calypso AI · Lakera · LASSO · NOMA · Pillar · AIM · ActiveFence · AIShield Powered by Bosch · Virtue AI · Acuvity · Deepkeep · Guardrails AI · Javelin · Arthur · Galileo · APEX · Tumeryk AI Trust Score · Haize Labs · pangea · Enkrypt AI · HIDDENLAYER · Liminal · Straiker · TROJ.AI · Operant
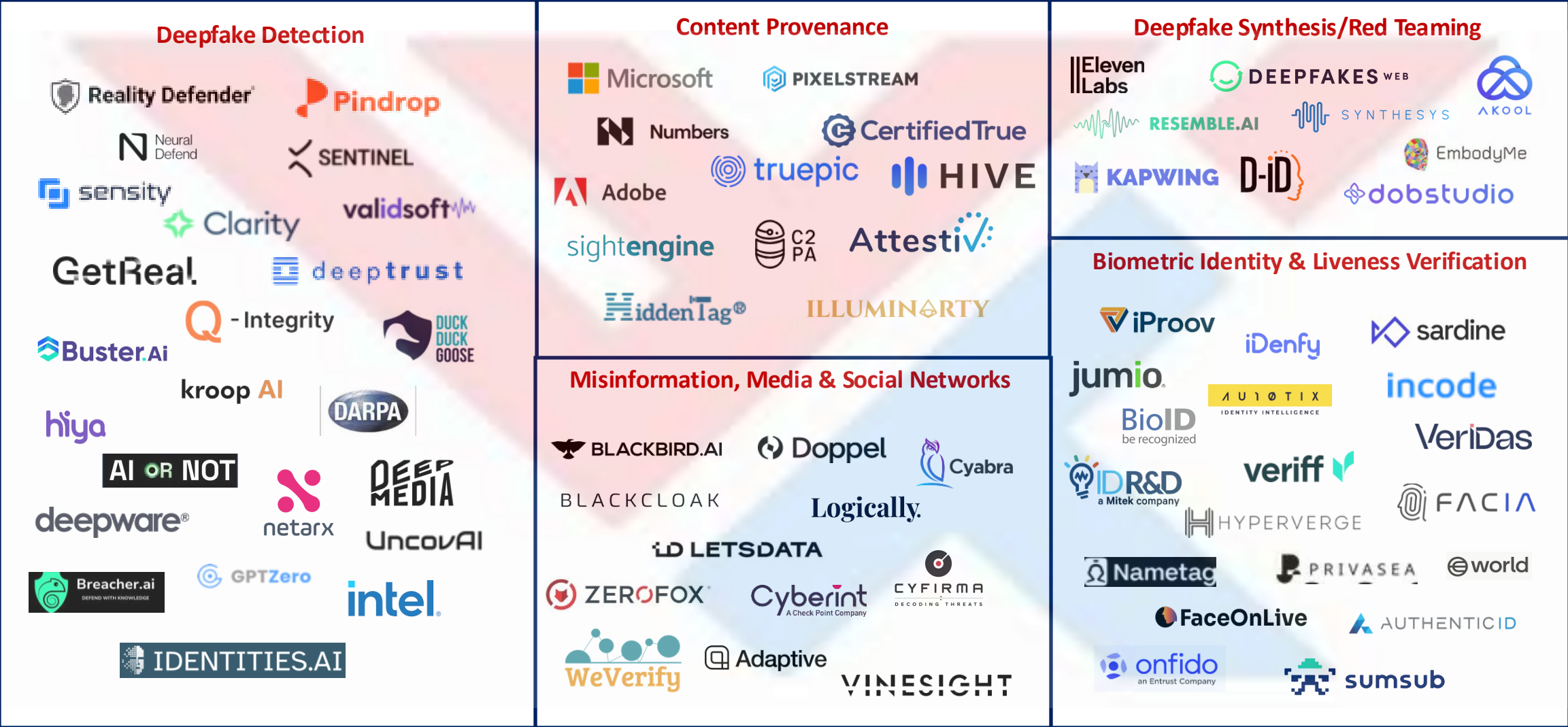
# AI for Cyber Security: WWT Market Landscape (v0512.2025)

# Security from Adversarial AI (Deepfake & Misinformation): WWT Market Landscape (v05.13.2025.76)

Thanks for Attending

**World Wide Technology**

**Todd Hathaway**

**Todd.Hathaway@wwt.com**

www.linkedin.com/in/toddhathaway