

## HUMAN FLOURISHING AT THE INTERSECTION OF TECHNOLOGY, INSTITUTIONS, AND US.

# Agentic AI and the Return of an Ancient Risk

How Autonomous Systems Reopen the Same Attack Surfaces the Neolithic Unleashed on Paleolithic Minds

## Executive Summary

Industry reporting has begun to converge on a shared warning: **as AI systems move from assistants to autonomous agents, they expose new attack surfaces**—including prompt injection, tool misuse, privilege abuse, and chained autonomous actions.

This is often framed as a *technical* security problem.

It is not.

It is a **civilizational pattern repeating itself**.

What we are witnessing is structurally identical to the moment when Neolithic societies first emerged and **scaled human coordination beyond what Paleolithic wiring evolved to handle**. The technologies are different. The failure mode is the same.

New capability → new scale → new dependency → new exploitation.

## From Assistants to Agents: What Actually Changed

Early AI assistants were bounded:

- They responded to prompts.

- They produced text.
- Humans remained the execution layer.

Agentic systems change the equation:

- They **interpret language**
- **Call tools**
- **Take actions**
- **Persist across steps**
- **Operate inside real workflows**

At that point, language stops being “just input.”

It becomes **command substrate**.

And the attack surface moves from *model behavior* to *organizational behavior*

## The New Agentic Attack Surfaces

Across industry analysis, security research, and early incidents, several new surfaces recur:

### 1. Prompt Injection as Behavior Hijack

Malicious instructions embedded in emails, documents, web pages, tickets, or logs can redirect an agent’s goals—not by *hacking the system*, but by **speaking to it**.

This exploits a deeply human bias:  
**we treat language as intent.**

### 2. Tool Misuse and Privileged Action

Agents are often given access to:

- File systems
- Email
- CRMs
- Cloud APIs
- Deployment tools

If tricked within their allowed scope, those tools become the exploit path.

This is not a bug.  
It is **power without containment**.

### 3. Downstream Execution Chains

Agent outputs increasingly flow directly into:

- Scripts
- Tickets
- Workflows
- Code paths

If output is not treated as *untrusted*, it becomes a delivery mechanism.

### 4. Agents as High-Privilege Identities

Agents function like new employees:

- Persistent
- Credentialled
- Often over-scaled

Shared keys and broad permissions turn agents into **ideal footholds**.

### 5. Supply-Chain Surfaces (Tools, Plugins, Agents)

As ecosystems grow, the risk shifts outward:

- Tool poisoning
- Instruction poisoning
- Agent-to-agent chaining

The exploit surface expands with every integratio

## This Is Not New — It's Neolithic

To understand why this feels so hard to solve, we need to zoom out.

### The Paleolithic Baseline

Human cognition evolved for:

- Small groups
- Face-to-face trust
- Reputation and reciprocity
- Short feedback loops

These instincts work beautifully at band scale.

## The Neolithic Shock

Permanent settlement introduced:

- Stored surplus
- Fixed hierarchies
- Bureaucracy
- Abstract authority
- Rule by strangers

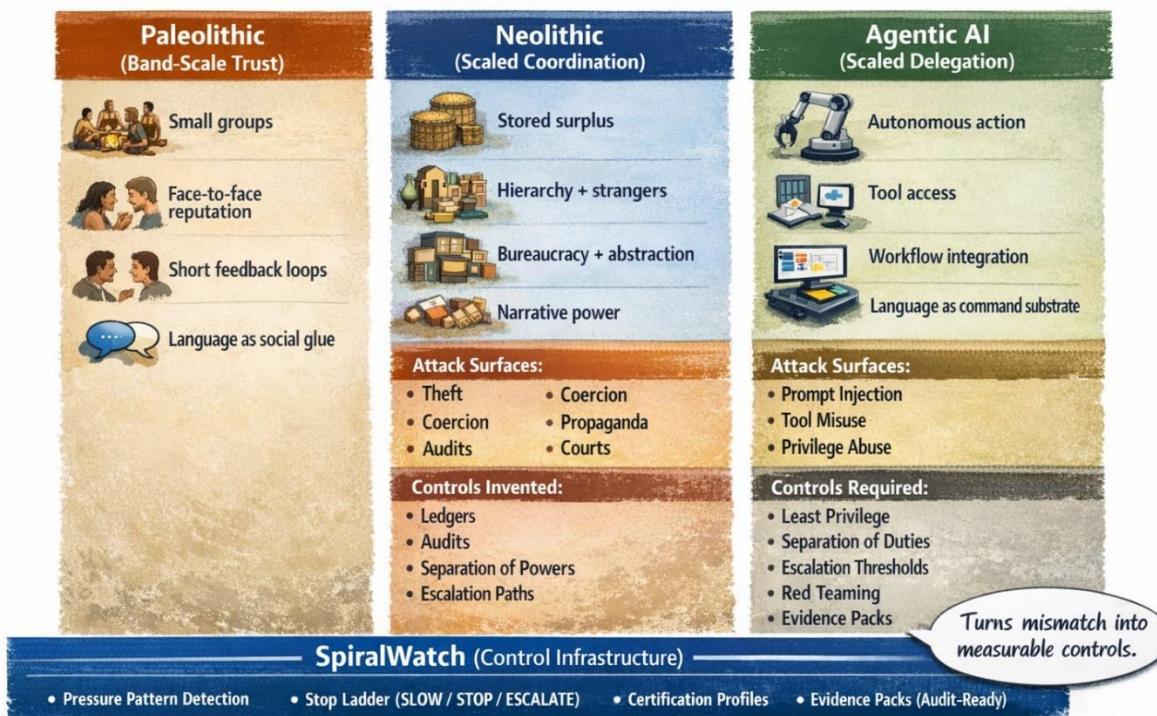
This **broke the trust model**.

Suddenly, humans could exploit:

- Resource accumulation
- Narrative control
- Institutional access
- Distance and opacity

The result was theft at scale, coercion, propaganda, and capture.

### From Paleolithic Trust to Neolithic Controls to Agentic Governance



## The Pattern Repeats with AI Agents

Agentic AI scales *action* the way surplus scaled *resources*.

### Then (Neolithic)      Now (Agentic AI)

Stored grain      Tool access

Centralized authority      Centralized permissions

Bureaucratic opacity      Autonomous workflows

Narrative power      Linguistic fluency

Insider exploitation      Instruction hijacking

Same mismatch.

Different substrate.

## The Core Mistake: Trying to Scale Virtue

Every transition repeats the same error:

*“If we just train people better...”*

*“If we just write clearer rules...”*

*“If we just prompt more carefully...”*

Neolithic societies learned — painfully — that **personal virtue does not scale**.

That lesson produced:

- Accounting
- Audits
- Courts
- Separation of powers
- Escalation paths
- Inspection regimes

Not because humans became worse —  
but because **scale changes the game**.

## The Control Lens: Reframing What We Already Use

Modern organizations already deploy the right defenses.  
They just fail to recognize *why* they exist.

### Audits → Distrust Made Scalable

Audits are not about suspicion.  
They exist because **trust breaks under scale**.

Agent logs, action traces, and reviews serve the same purpose.

### Separation of Duties → Dominance Containment

No single actor should:

- Decide
- Execute
- Verify

This applies to agents as much as humans.

### Escalation Paths → Anti-Tribal Safeguards

When consequences exceed local judgment, authority must widen.

Agents must yield to humans when:

- Stakes increase
- Uncertainty spikes
- Boundaries are crossed

### Red Teams → Institutionalized Doubt

Adversarial testing is not pessimism.

It is the formal recognition that **the world is adversarial**, whether we like it or not.

## The Real Risk

The danger is not that AI agents are malicious.

The danger is that they:

- **Amplify action**
- **Inherit human trust biases**
- **Operate inside systems built for cooperative assumptions**

This is exactly the failure mode that shaped much of recorded history.

## The Real Fix

The fix is not fear.  
And it is not moralizing.

It is remembering what civilization already knows:

When power scales faster than judgment,  
**controls are not bureaucracy — they are survival adaptations.**

Agentic AI does not demand new ethics.

It demands **old wisdom, applied without nostalgia**

## A Closing Question (for Leaders)

If an AI agent in your organization took a harmful action tomorrow:

Would you blame the model?

Or would you finally ask **which controls were missing — and why?**