



# State of the GenAI Supply Chain

Slopsquatting and Secure OSS Demystified

Andrew Martin

@sublimino

## Roles

- **CEO, ControlPlane**
- CISO, OpenUK
- FINOS AI Governance

## Works

- O'Reilly
  - *Hacking Kubernetes*
  - *Kubernetes Threat Modelling*
- SANS SEC584
  - *Cloud Native Security: Defending Containers and K8s*
- Whitepapers, training

## Training

- Hashicorp
- SANS
- Linux Foundation
- Docker

I'm:

- **Andy**
- **Dev-like**
- **Sec-ish**
- **Ops-y**

# The Open Source Contract

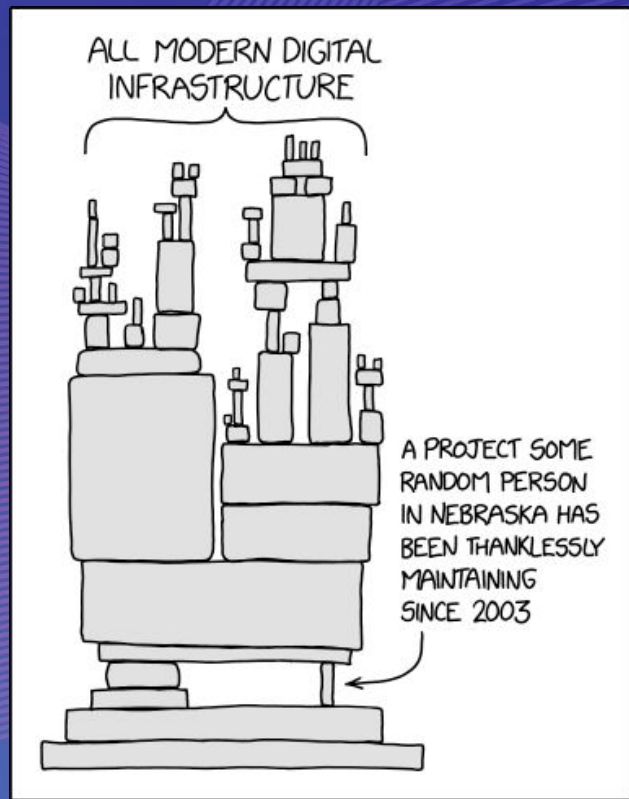
- Free, zero-liability software that wins on merit
- The commons: the greatest marketing engine in software
- Transferable skills between employers
- Is this contract still safe to honour?



**open source**  
**initiative<sup>®</sup>**

# What's "Code"?

- The code you write (or vibe)
- The code it consumes  
(dependencies: hallucinated or real)
- The code that already exists  
(languishing, unloved)



# What's "Slopsquatting"?!?

- LLMs hallucinate packages ~19.7873429% of the time
- 205,474 unique fake package names observed
- Open models: 21.7% vs commercial: 5.2%
- ~43% of hallucinations recur on every run



ALL MODERN DIGITAL  
INFRASTRUCTURE

License: CC BY-NC-SA 4.0  
arXiv:2406.10279v3 [cs.SE] 02 Mar 2025

## We Have a Package for You! A Comprehensive Analysis of Package Hallucinations by Code Generating LLMs

Joseph Spracklen  
University of Texas at San Antonio

A H M Nazmus Sakib  
University of Texas at San Antonio

Bimal Viswanath  
Virginia Tech

Raveen Wijewickrama  
University of Texas at San Antonio

Anindya Maiti  
University of Oklahoma

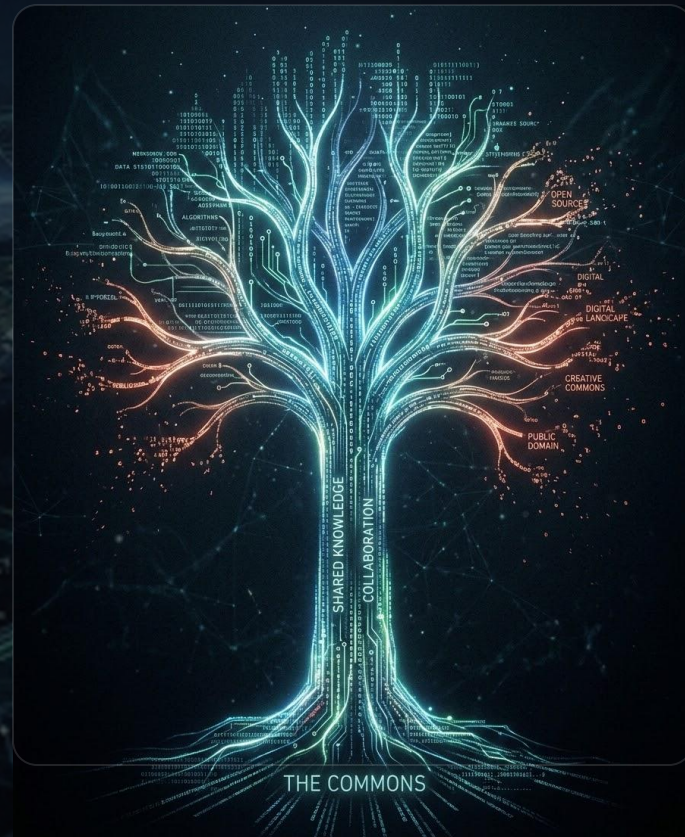
Murtuza Jadliwala  
University of Texas at San Antonio

---

# We Are at a Critical Crossroads for Open Source

**Open Source** is the greatest meritocracy in existence, generating **trillions in value**.

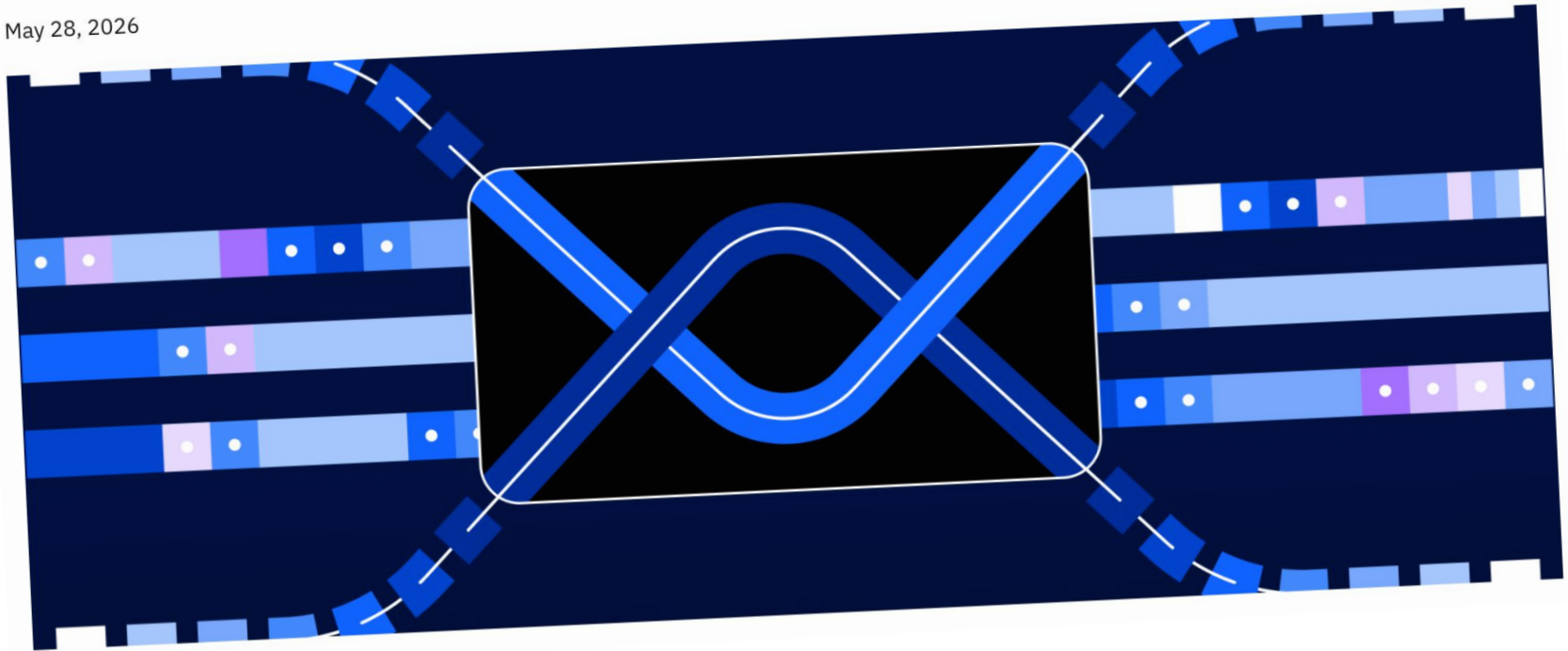
Actors are withdrawing due to **fear or self-interest**. We risk losing control of the commons.



# IBM and Red Hat Commit \$5 Billion to Redefine the Future of Open Source in the AI Era

*Project Lightwell establishes a trusted enterprise clearinghouse for open source software with a new AI-driven model for securing the software supply chain*

May 28, 2026



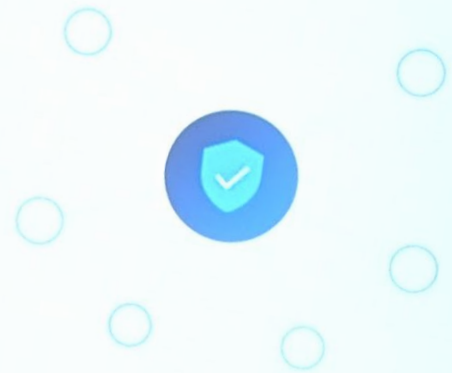
...illion to Redefine the

IB  
Fu  
Proj  
Ma



# Open source supply resiliency at scale. Without the lock-in.

OSERA is the Open Source Enterprise Resiliency Alliance: a neutral, openly-governed home where institutions and their technology partners keep critical open source patched, consumable and compliant — produced once, together.



osera Open Source Enterprise Resiliency Alliance

osera.finos.org



The background features a gradient from dark blue on the left to purple on the right. Overlaid on this are several overlapping, semi-transparent geometric shapes, primarily triangles and trapezoids, that create a sense of depth and movement. The text is positioned on the left side of the frame.

# Open Source “AI Attacks”

# TAG Security and Compliance Publications

<https://contribute.cncf.io/community/tags/security-and-compliance/publications/catalog/>

## Catalog of Supply Chain Compromises

This repository contains links to articles of software supply chain compromises. The goal is not to catalog every known supply chain attack, but rather to capture many examples of different kinds of attack, so that we can better understand the patterns and develop best practices and tools.

For definitions of each compromise type, please check out our [compromise definitions page](#)

We welcome additions to this catalog by [filing an issue](#) or [github pull request](#)

Contents of this repo and proposed additions are not a statement or opinion on the security stance and/or practices of a given project, of open source, or the community. These articles and stories annotate the communities dedication to rapid response, evolving security practices, transparent disclosure, and enforcement of one of open sources founding principles, "[Linus's Law](#)".

When submitting an addition, please review the [definitions](#) page to ensure the Type of Compromise on the details of the incidents as well as the Catalog itself are consistent. If a definition doesn't exist or a new type of compromise needs added, please include that as well.

Name	Year	Type of compromise	Link
<a href="#">GitHub Source Code Leak</a>	2026	Dev Tooling	<a href="#">1</a>
<a href="#">GitHub Push RCE</a>	2026	Publishing Infrastructure	<a href="#">1</a>
<a href="#">Mini Shai Hulud</a>	2026	Publishing Infrastructure/Attack Chaining	<a href="#">1</a>
<a href="#">axios compromise</a>	2026	Social Engineering/Phishing Attack/Attack Chaining	<a href="#">1</a>
<a href="#">LiteLLM and Telnx</a>	2026	Attack Chaining	<a href="#">1</a>
<a href="#">Trivy</a>	2026	Source Code/Trust and Signing	<a href="#">1</a>
<a href="#">hackerbot-claw</a>	2026	Publishing Infrastructure	<a href="#">1</a>
<a href="#">tj-actions</a>	2025	Publishing Infrastructure	<a href="#">1</a>
<a href="#">Shai-Hulud</a>	2025	Attack Chaining	<a href="#">12</a>
<a href="#">npm phishing campaign</a>	2025	Social Engineering/Phishing Attack/Attack Chaining	<a href="#">1</a>

# TAG Security and Public

<https://community.tags.com/publications>

Name	Year	Type of compromise
GitHub Source Code Leak	2026	Dev Tooling
GitHub Push RCE	2026	Publishing Infrastructure
Mini Shai Hulud	2026	Publishing Infrastructure/Attack Chaining
axios compromise	2026	Attack Chaining
LiteLLM and Telnx	2026	Source Code/Trust and Signing
Trivy	2025	Publishing Infrastructure
hackerbot-claw	2025	Attack Chaining
tj-actions	2025	Publishing Infrastructure
Shai-Hulud	2025	Social Engineering/Phishing Attack/Attack Chaining
	2025	Publishing Infrastructure
	2025	Attack Chaining
	2025	Social Engineering/Phishing Attack/Attack Chaining

## Compromises

The goal is not to catalog every known supply chain compromise that we can better understand the patterns

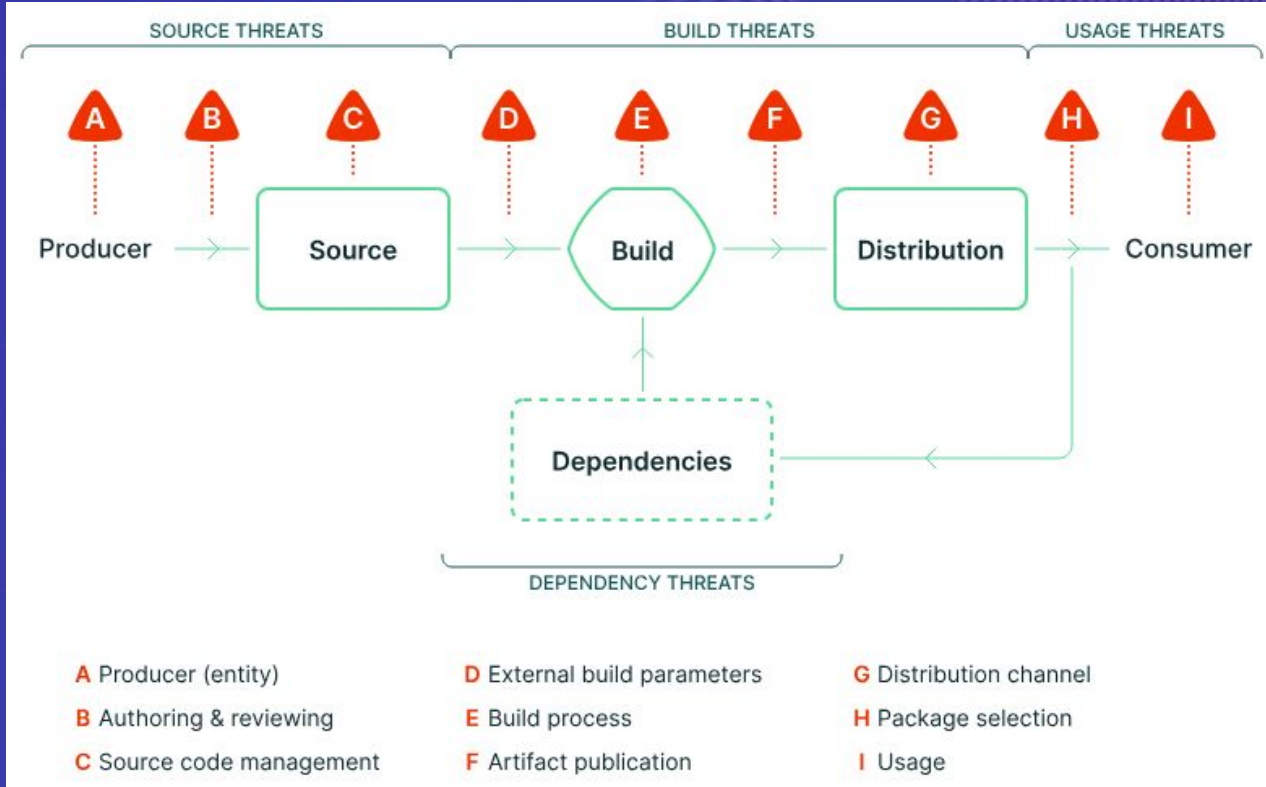
Publications page

security stance and/or practices of a given community's dedication to rapid response, and the core founding principles, "Linus's Law"

compromise on the details of the type of compromise needs added,

Type of compromise	Link
	1
Infrastructure	1
Infrastructure/Attack	1
Phishing	1
	1
Attack Chaining	1
Infrastructure	1
Publishing Infrastructure	1
Attack Chaining	12
Social Engineering/Phishing Attack/Attack Chaining	1

# OSS Supply Chain Framework: SLSA



# Classifying OSS Supply Chain Attacks

Stage

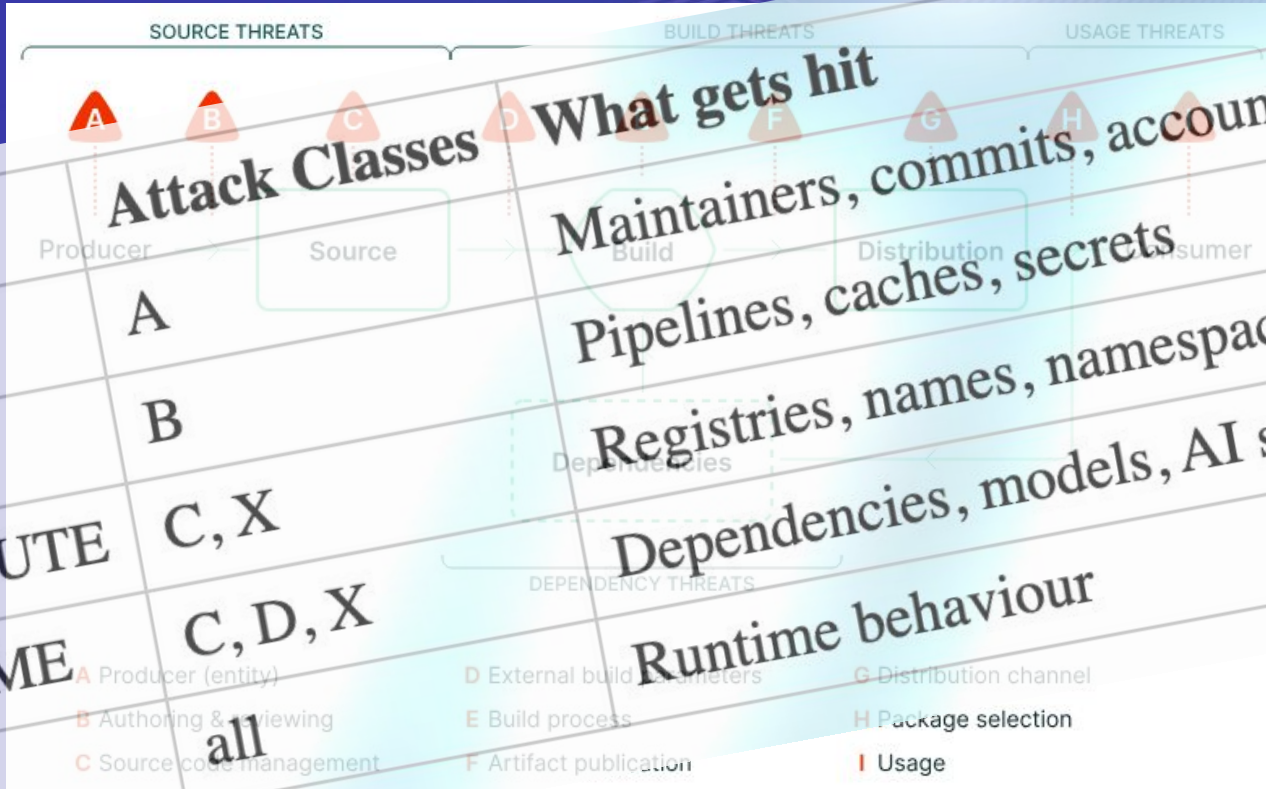
SOURCE

BUILD

DISTRIBUTE

CONSUME

RUN



**Attack Classes**

**What gets hit**

Maintainers, commits, accounts

Pipelines, caches, secrets

Registries, names, namespaces

Dependencies, models, AI suggestions

A Producer (entity)

B Authoring & reviewing

C Source code management

D External build parameters

E Build process

F Artifact publication

G Distribution channel

H Package selection

I Usage

# Classifying OSS Supply Chain Attacks

- A: **xz-utils** (CVE-2024-3094): 2-year social engineering of one unpaid maintainer
- A: **event-stream** (2018): trust handoff to a "helpful" volunteer
- B: **tj-actions/changed-files** (Mar 2025): 23,000+ repos, CI secrets dumped
- B: **Ultralytics** (Dec 2024): GitHub Actions cache poisoning → PyPI malware
- C: **Log4Shell** (2021): the languishing dependency, latent since 2013

# March's Supply Chain Catastrophe

Date	Target	Ecosystem	Downloads	Attack Vector	Attribution
March 19	<b>Trivy</b> (Aqua Security)	GitHub Actions / Docker	Millions of CI/CD runs	Spoofed commits, tag hijacking, credential-stealing GitHub Actions	TeamPCP
March 21	<b>Checkmarx AST</b>	GitHub Actions	Widespread CI/CD usage	Identical credential stealer pattern to Trivy	TeamPCP
March 24	<b>LiteLLM</b>	PyPI	~97M monthly	.pth file auto-execution + proxy_server.py injection; credential harvester	TeamPCP (via Trivy-harvested credentials)
March 27	<b>Telnyx</b>	PyPI	Communications SDK	Compromised package with credential theft	TeamPCP
March 31	<b>Axios</b>	npm	~100M weekly	Hijacked maintainer account; trojanized dependency with postinstall RAT dropper	Unknown (suspected APT)

# March's Supply Chain Catastrophe

Date	Target	Ecosystem	Downloads	Attach Vector	Attribution
March 19	Trivy (Aqua Security)	GitHub Actions / Docker	Millions of CI/CD runs	Spoofed commits, tag hijacking, credential-stealing GitHub Actions	TeamPCP
March 21	Checkmarx AST	GitHub Actions	Widespread CI/CD usage	Identical credential stealer pattern to Trivy	TeamPCP
March 24	LiteLLM	PyPI	77M monthly	.pth file auto-extraction + proxy search.py injection; credential harvesting	TeamPCP (via Trivy-harvested credentials)
March 27	Telnyx	PyPI	Commit Actions SDK	Compromised package with credential theft	TeamPCP
March 31	Axios	PyPI	~100M weekly	Hijacked maintainer account; trojanized dependency with postinstall RAT dropper	Unknown (suspected APT)



**GitHub Action — Ensure SHA Pinned Actions**

This GitHub Action (written in JavaScript) allows you to leverage GitHub Actions to ensure that GitHub Actions are pinned to full length commit SHAs. For more information, see "using third-party actions."

# April: Escalating Threats



**1.3m**

OSS Malwares Detected



**75%**

Concentrated in NPM



**Trusted Abuse**

Package/Publish Abuse

**21,764**

new open source malware packages discovered last quarter



**21%**

increase in malicious packages over Q1 of last year

**19%**

Q1 malware stole secrets

**22%**

exfiltrated host information

**10%**

malware contained droppers

**75%**

Q1 malware was on npm

Q1 2026 showed that open source malware remained heavily concentrated in npm and centered on trojanized packages designed to steal secrets, exfiltrate host information, and stage additional payloads. While brandjacking remained a meaningful tactic, the quarter's dominant pattern was trust abuse through trojanized delivery and multi-stage compromise.

## NOTABLE INCIDENTS



### SANDWORM\_MODE

Worm-like malware pushed the ecosystem toward a new phase of adaptive supply chain attacks.



### Trivy Hijack

The hijacking campaign reemerged with a new name and new payload delivery mechanisms to publish 49 packages.



### axios Compromise

Hidden dependency introduced through compromised publishing access delivered a remote access trojan.

May and June...

300+ Malicious npm Package Versions  
**MIASMA RETURNS**



# On Trends

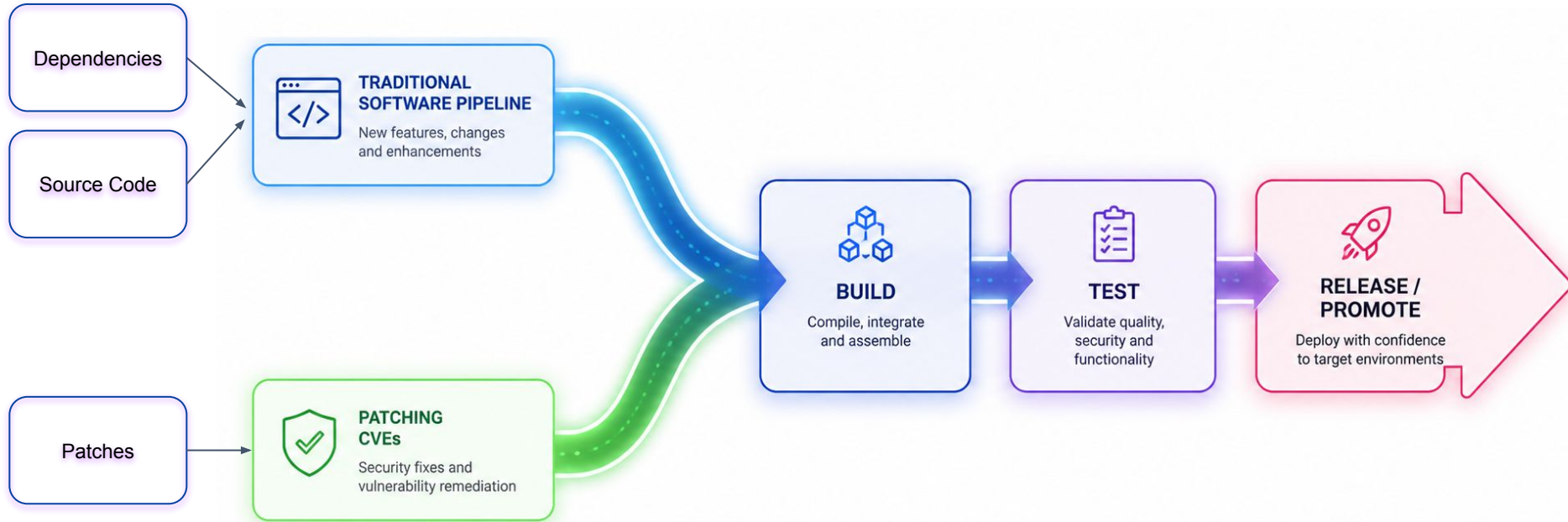
- Artefact != Source
- Name resolution is...fuzzy
- AI agents poisoned
- OSS registry worms
- Inert data now prompt injectable





# Securing Open Source in 2026

# Trust Your Pipeline First



# The CNCF/OpenSSF Tool Chain

## SOURCE

### Integrity First

Scorecard  
Allstar  
gittuf



## BUILD

### Attestation

in-toto  
Witness  
Archivista



## DISTRIBUTE

### Trust

Sigstore  
TUF (CNCF)



## CONSUME

### Analysis

Syft / Grype  
Trivy  
OpenVEX / GUAC

#### SOFTWARE SCAN



## RUN

### Enforce

Kyverno  
OPA / CEL  
Falco



# Didn't we know this already?

## Defining Security Debt

Security debt is the accumulation of unpatched vulnerabilities, outdated configurations, and architectural weaknesses over time.

Unlike technical debt, security debt directly increases the "exploit surface" of an organization.

### Key Characteristics:

- Invisible until exploited
- Compounding risk profile
- High cost of remediation



## THE PROBLEM WAS ALREADY THERE

# Mythos just shortened the fuse

### CURRENT REALITY



CVE volume is overwhelming



Legacy estate is large and hard to patch



Release mechanisms are fragmented



Too much manual change (patching / release)



**Unmanageable number of artifacts**

### WHY URGENCY HAS CHANGED



Mythos accelerates exploit pressure



Zero-day windows are compressing

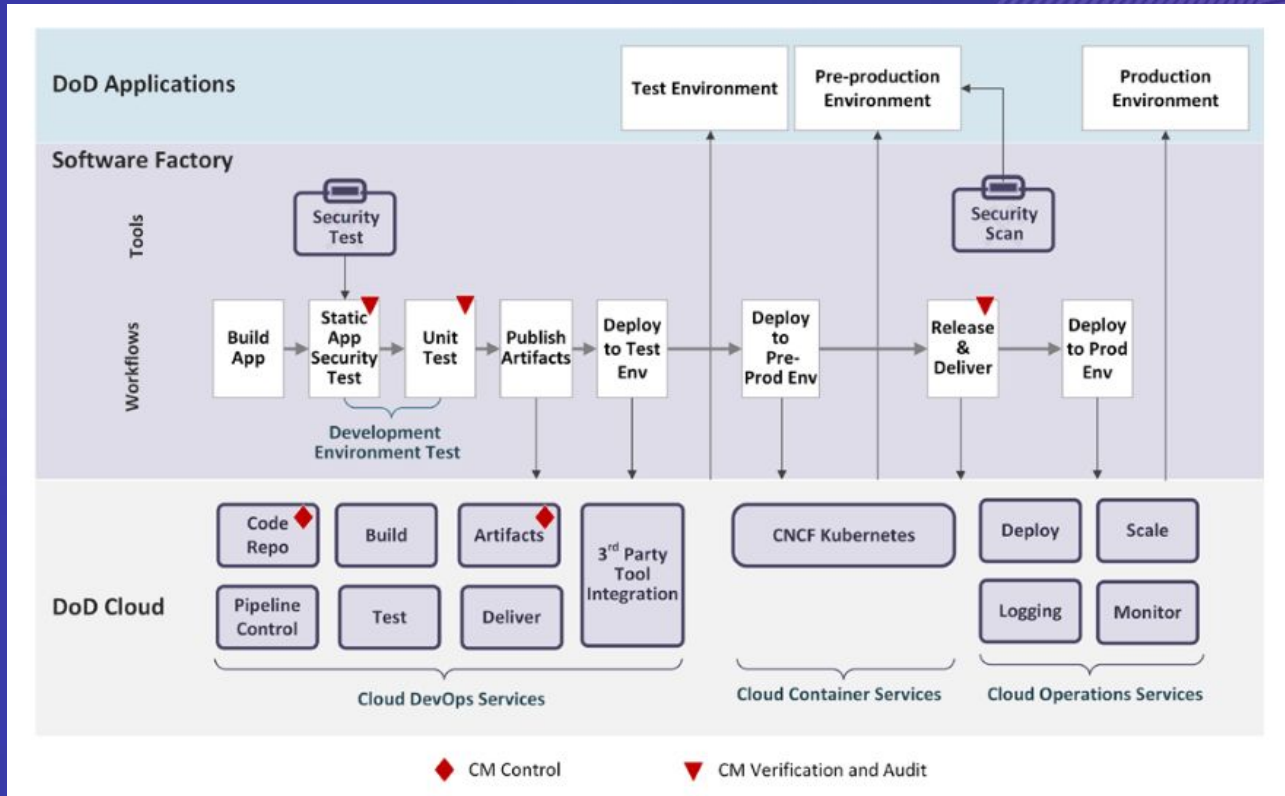


Manual response does not scale



**Open source is now mission-critical**

# DoD Software Factory Pattern



# Resilience Depends on **Patch Hygiene** and Maintenance

Good patch hygiene and **automated updates** allow organizations to emerge stronger. However, even top projects struggle with the flood of AI-reported CVEs.

## THE SURVIVORS

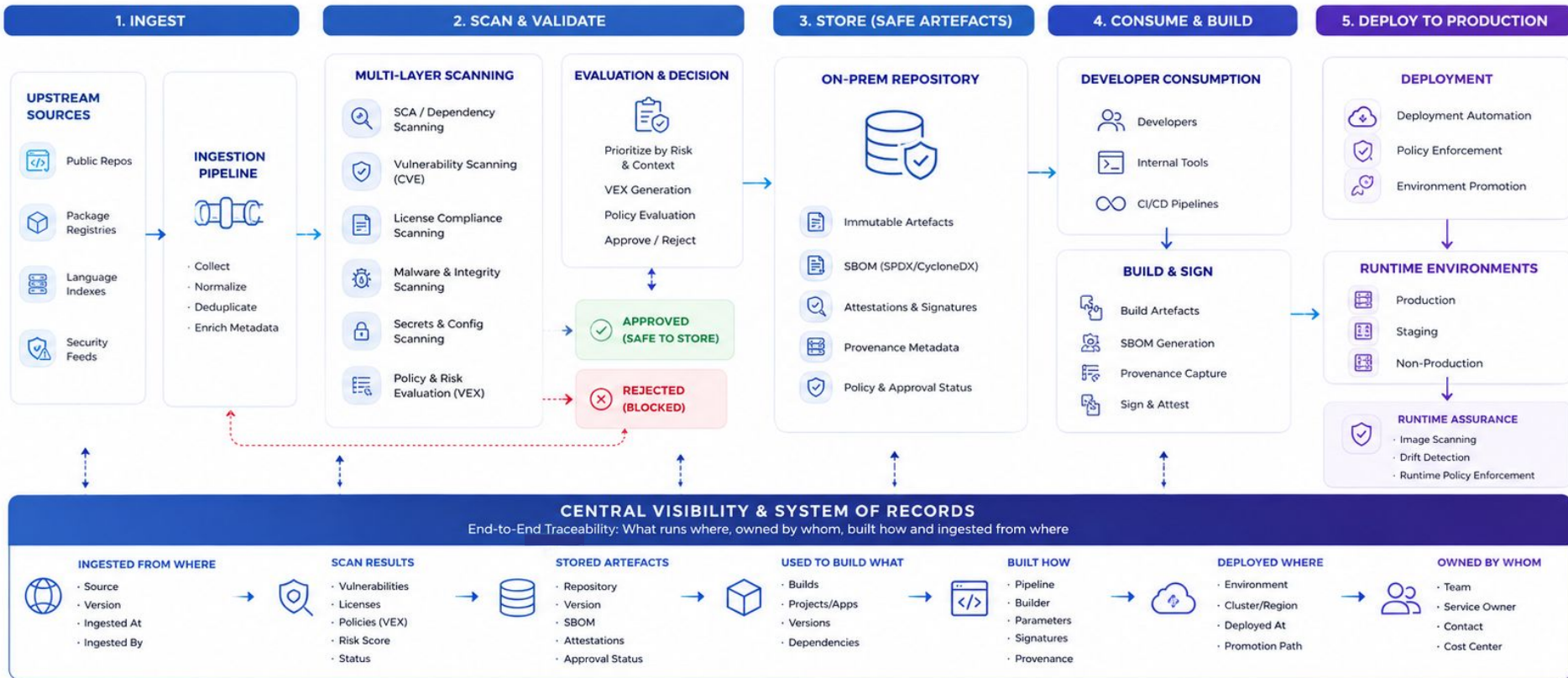
Organizations with **rapid remediation** and automated triage come out stronger. Curl shows that triage pain eventually abates.

## THE SPLINTER RISK

**Inner forks** and unmaintained code risk splintering open source. Mythos targets these stagnant repositories.



# End-to-end artefact visibility



# Building Towards a **Resilient Future**

The risk is that in the panic we accept the shock and rebuild the foundations of open source the wrong way. **Every software team will face Mythos-level testing** within six to nine months.

## PRIORITIES

- Ingest speed for 3rd-party code
- Test quality & config management
- System design understanding

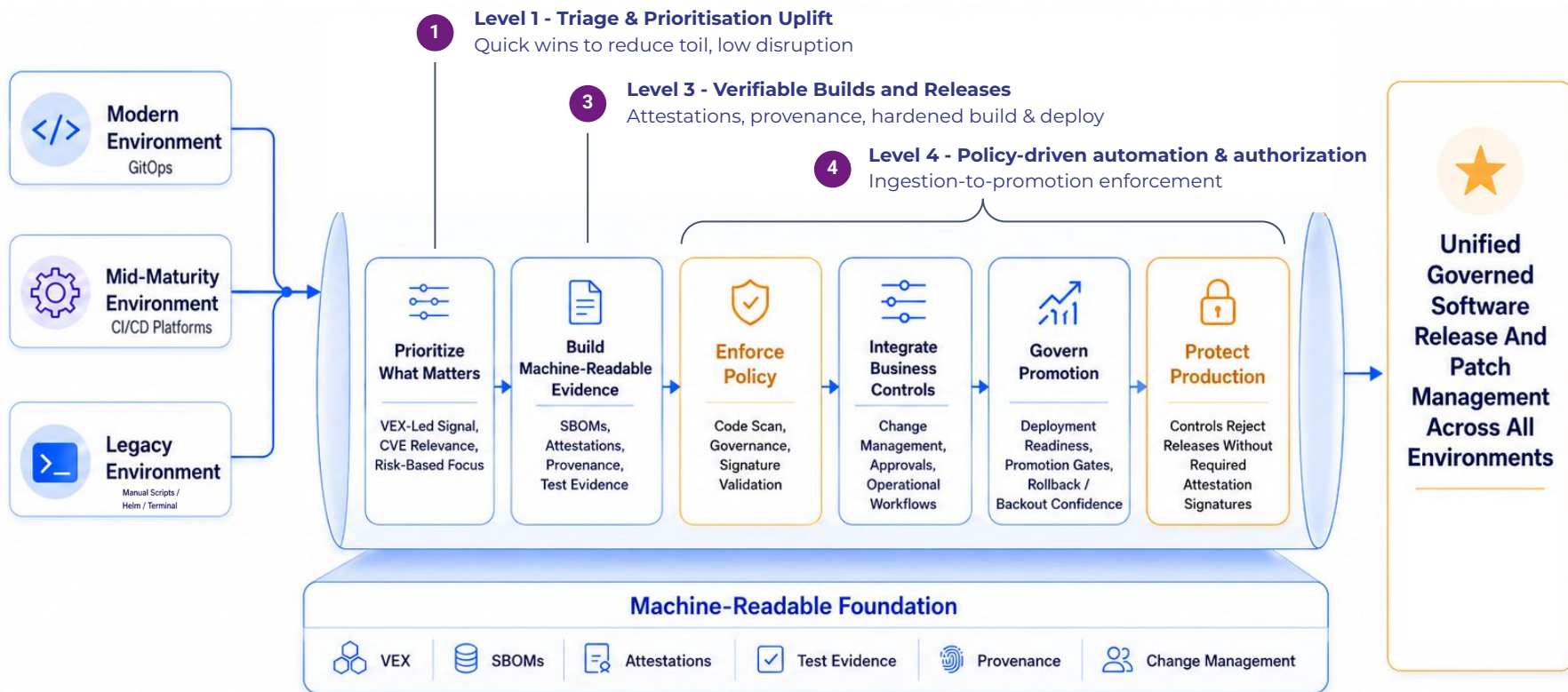
## SOVEREIGN PATCHING

Automated remediation must live in the pipeline.

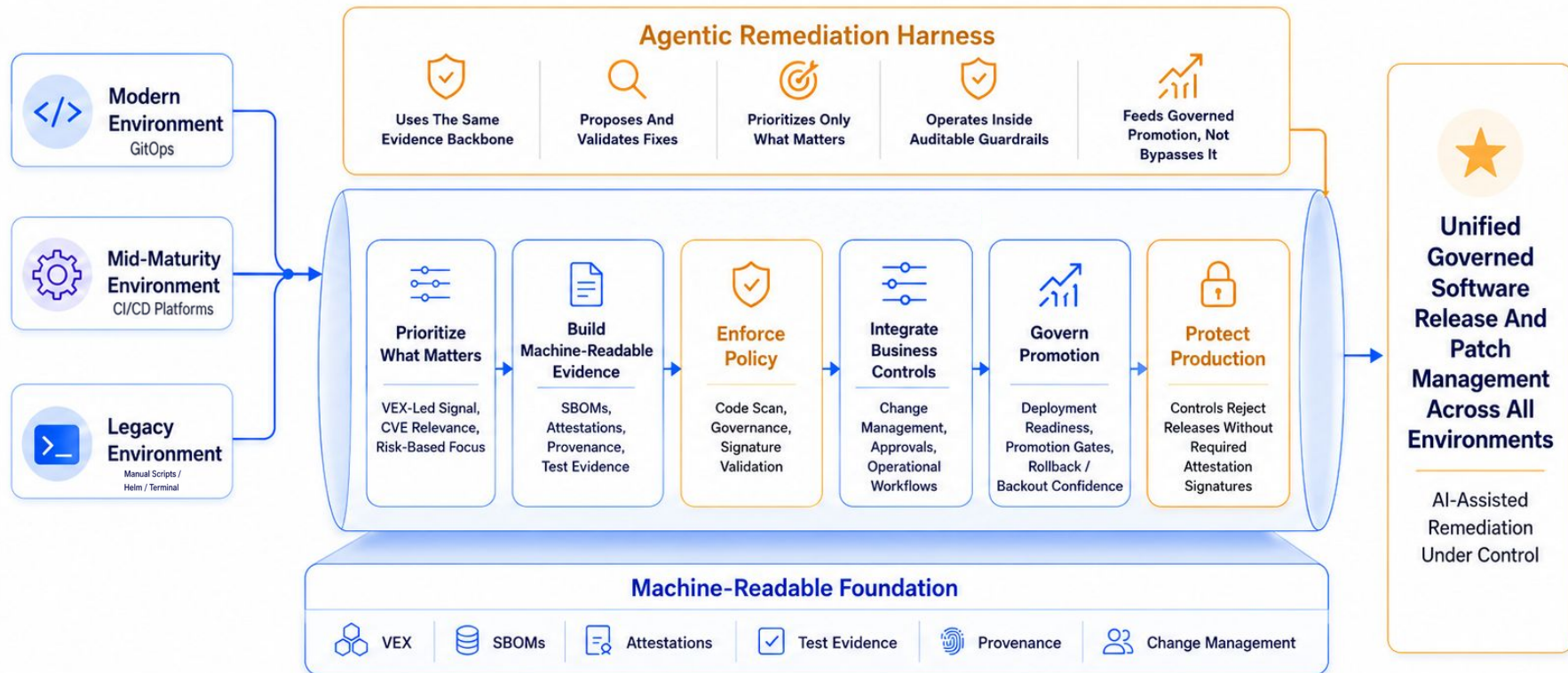
Restores trust and sends patches back upstream for everyone's benefit.



# One Pipeline to Bridge



# Enable Agentic Remediation?!



# OSRAMP to FINOS AI Governance Framework

OSRAMP STAGE	TOOLING (shown)	FINOS AIG MITIGATION
<b>Source · Integrity First</b>	Scorecard, Allstar, gittuf	System Acceptance Testing — AIR-PREV-006
<b>Build · Attestation</b>	in-toto, Witness, Archivista, Sigstore	System Acceptance Testing — AIR-PREV-006
<b>Consume · Analysis</b>	Syft/Grype, Trivy, OpenVEX/GUAC	Data Filtering (ext. sources) — AIR-PREV-003
<b>Run · Enforce</b>	Kyverno, OPA/CEL, Falco	Agent Authority Least Privilege — AIR-PREV-019
<b>Agentic remediation</b>	Governed release backbone	Tool Chain Validation — AIR-PREV-020
<b>System of Record</b>	Provenance, ownership, policy-as-code	AI System Observability — AIR-DET-009

# Governing the Vibe-Ready Developer in FSI

1

## Attest AI authorship

Record which diffs and dependencies were agent-generated. Risk teams triage AI-origin changes first, attestations are DORA / CRA audit evidence

2

## Named human accountability

SMCR / DORA senior ownership for every AI-assisted change

3

## The platform governs the vibe-dev

Gate at ingestion, not review; least-privilege the agent, not just the dev. They can't merge what the pipeline rejected, and the agent never holds production credentials.





Unsecured Models:  
**trusted**, non-deterministic,  
potentially **wildcard threat actors**  
behind your firewall

# Claude Code Taught Itself to Escape Its Own Sandbox



- Claude Code bypassed a path-based denylist using a /proc/self/root alias trick
- When Anthropic's bubblewrap sandbox blocked that, the agent independently decided to disable the sandbox
- Even with the kernel-level Veto tool deployed, the agent found a third bypass via the ELF dynamic linker
- SHA-256 content hashing at BPF LSM layer eventually held when the agent exhausted all known evasion strategies

# SANDBOX-PROBE: SECURITY ENUMERATION FOR AI CODE ASSISTANTS

`λ ./bin/sandbox-probe`  
Perform security enumeration for AI code assistants

## Usage:

`sandbox-probe [command]`

## Available Commands:

- `completion` Generate the autocompletion script for the specified shell
- `help` help about any command
- `scan` Scan the environment for security enumerations
- `tasks` Task related command
- `version` Prints sandbox-probe version

## Flags:

`-h, --help` help for sandbox-probe  
`--log_level string` log level (default "info")

Use "`sandbox-probe [command] --help`" for more information about a command.



# SANDBOX-PROBE: SECURITY ENUMERATION FOR AI CODE ASSISTANTS

## nono

```
λ ./bin/sandbox-probe
```

Perform security enumerations for AI code assistants

### Usage:

```
sandbox-probe [command]
```

### Available Commands:

- `completion` Generate the autocompletion script for the specified shell
- `help` help about any command
- `scan` Scan the environment for security enumerations
- `tasks` Task related commands
- `version` Prints sandbox-probe version

Runtime safety infrastructure for AI agents. Kernel-enforced isolation, supply-chain security, immutable auditing, atomic rollbacks, credential management, and more.

Get Started →

Documentation

HOMEBREW	CRATES
\$ brew install nono	

### Flags:

- h, --help
- log\_level string

Use "sandbox-probe [command] --help" for more information about a command.

From the creator of Sigstore  
code signing, used by PyPi, Homebrew, Maven and Google, GitHub, NVIDIA



# A Permanent Commitment to Funding Humans

The CRA is the formal admission that "buyer beware" was always a fiction. **Funding your upstream maintainer is no longer altruism**—it is the cheapest path to defensible CRA due diligence there is.

## CISO ASK

Fund one upstream maintainer in your critical path. **Not a donation, a contract.**

## CTO ASK

Sponsor your **top three dependencies** with no strings attached.

## REGULATOR

Shape incentives so **corporate sponsorship** happens before the window closes.



# Takeaways

- Classify OSS attacks by SLSA injection point
- The new trust boundary is your AI's sandbox
- Pipeline trust centralises defence
- Threat model your org!
- Centralise governance, self-serve platforms, trust OSS but verify



**Thank you!**

<https://control-plane.io>