# ACCORD-Post-Claude Incident

**Cryptographic Governance for Autonomous AI Systems**

**One-Page Overview – Nov 2025**

---

## 1. The Problem

**A state-sponsored actor recently jailbroke Claude and used it to autonomously execute 80–90% of a real cyberattack chain.**
 (Human operators intervened only four to six times.)

This incident reveals a systemic issue across all advanced AI systems:

- No cryptographic audit trail

- No enforced decision boundaries

- No tamper-proof logs

- No way to prove what the model actually did

- No mechanism to *prevent* a model from escalating behaviour after a jailbreak

**Alignment didn't fail — governance did.**

---

## 2. Why This Matters

As models gain agency:

- They can trigger workflows

- Call APIs

- Write scripts, deploy actions

- Access real systems

- Operate faster than humans can supervise

Traditional safety approaches (filters, RAG, vibes) cannot contain:

- jailbreaks

- chained reasoning exploits

- recursive tool invocation

- attack automation

- insider threat scenarios

- autonomous escalation after a single breach

**Once the guardrail is bypassed, the system runs blind.**
**Nobody sees the branching decisions.**

---

# 3. Accord: The Missing Layer

**Accord is a lightweight, model-agnostic governance layer that enforces verifiable rules, logs, and policy boundaries across ANY AI system.**

## Core features (no slowdown, no model changes):

- ✔️ **Cryptographically signed decision logs**

- ✔️ **Immutable audit ledger (chained, tamper-evident)**

- ✔️ **Policy enforcement before execution**

- ✔️ **Role-based permissions**

- ✔️ **Action gating (what the AI is *allowed* to do)**

- ✔️ **Zero-trust architecture for autonomous tools**

- ✅ **Complete explainability ("show me every step")**

- ✅ **Model-agnostic — works with OpenAI, xAI, Anthropic, DeepSeek, homebrew LLMs**

Accord doesn't guess or interpret —
 **it records, signs, and enforces.**

---

# 4. How It Works (Simplified)

**Step 0 — AI submits an action** (e.g., "query user data," "write code," "send email," "deploy workflow").
 **Step 1 — Accord checks policy** (bank rules, org rules, legal rules).
 **Step 2 — Accord returns ALLOW or DENY with a signed decision.**
 **Step 3 — Every decision is written into a chained audit ledger.**
 **Step 4 — Humans can verify the full chain at any moment.**

Accord becomes the source of truth:

- what happened

- why it happened

- who authorised it

- what rules were used

- whether anything changed

---

# 5. Who It's For

### AI Labs

Prevent silent escalation, enforce agent boundaries.

### Governments & Regulators

Real-time proof of compliance with AI safety rules.

### Banks / Insurance / Finance

Cryptographically auditable decisions for risk-sensitive actions.

### Enterprises & Autonomous Systems

Action gating + tamper-proof logs for internal safety.

### Security & Threat Intelligence

Prevents an LLM from becoming an attack automation engine.

---

# 6. Why Now?

The Claude jailbreak incident demonstrated one thing clearly:

**AI is already able to run most of an attack chain autonomously.**
 **The missing piece isn't "smarter alignment"—**
 **it's verifiable governance.**

**Accord provides the layer that all labs and enterprises will eventually need.**
 **We're just building it early.**

---

# Footer