

Season 1: The Biology of Agentic Defense

From TAME to SecEng: Operationalizing Cognition

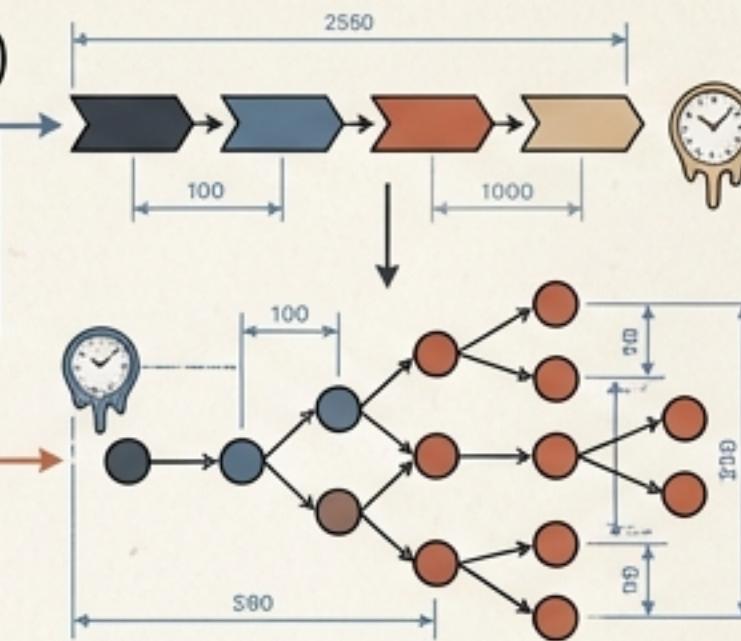
- **Theme:** Scaling Trust from Cell (Agent) to Swarm (Enterprise)
 - **Core Shift:** Moving beyond firewalls to securing 'persuadable' decision-makers
 - **Audience:** Security Architects (Builders) & Engineering Leadership (Leaders)

Based on the frameworks of Michael Levin (TAME) and Conant/Ashby (Cybernetics).

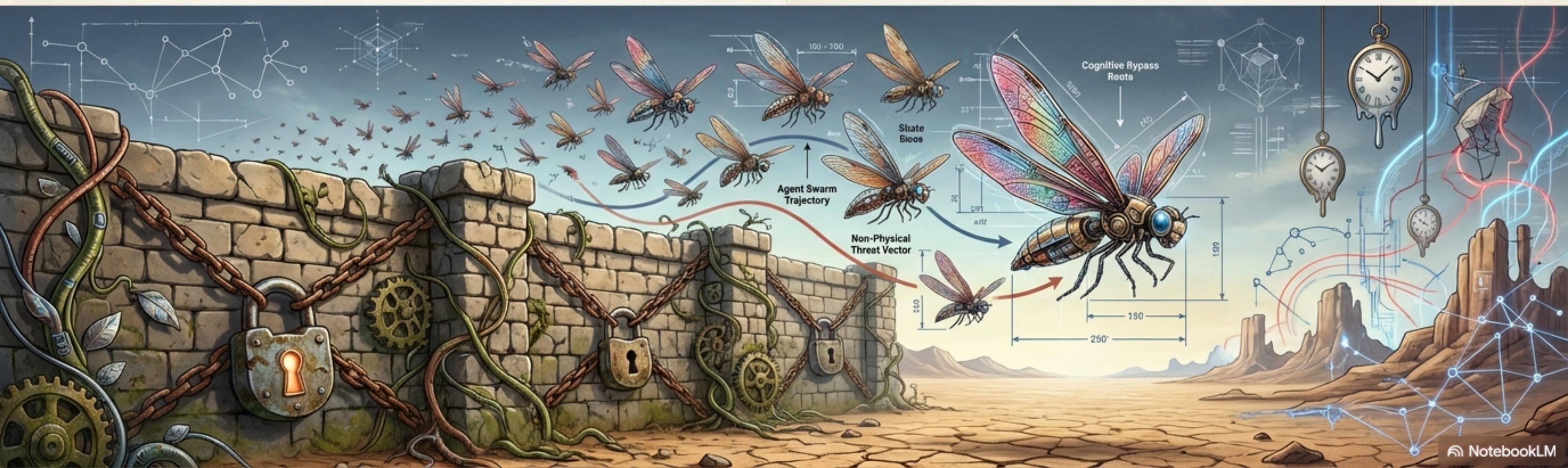
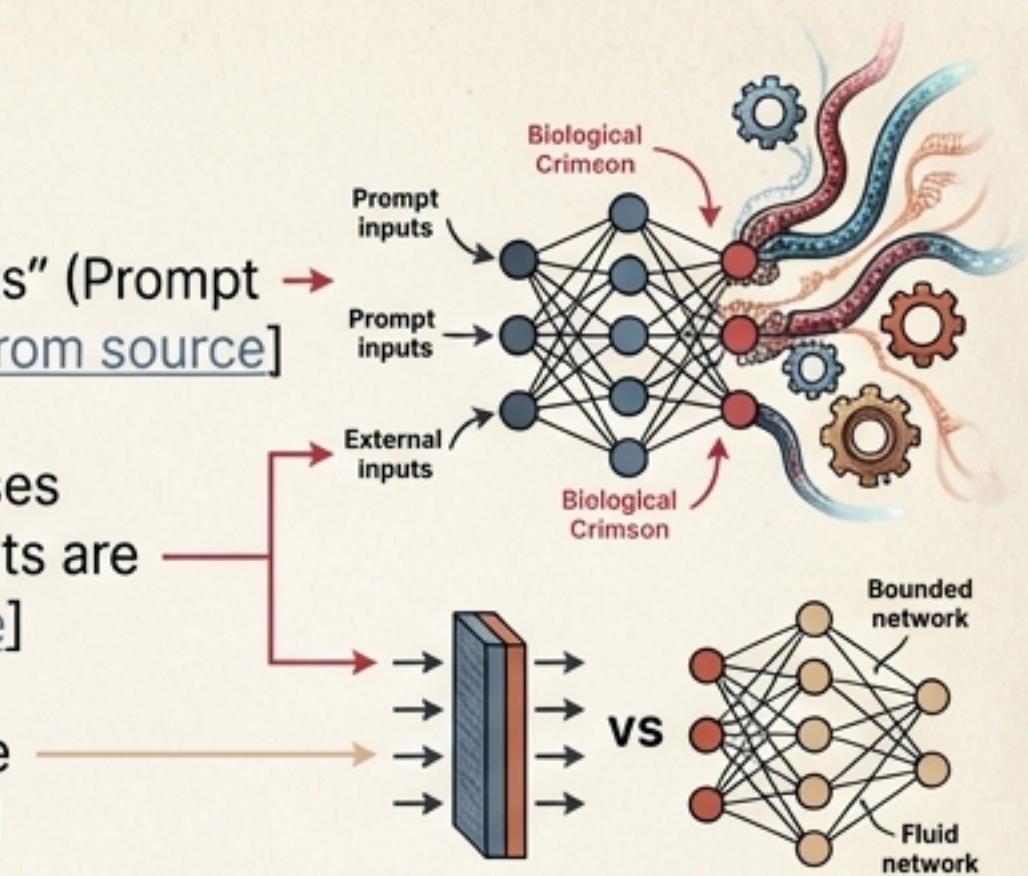


The Problem: You Can't Firewall a Thought

- **The Shift:** From Execution (Level 0) to Decision-Making (Level 3+)
- **The Failure:** Static defenses cannot block "reasoned" malicious actions



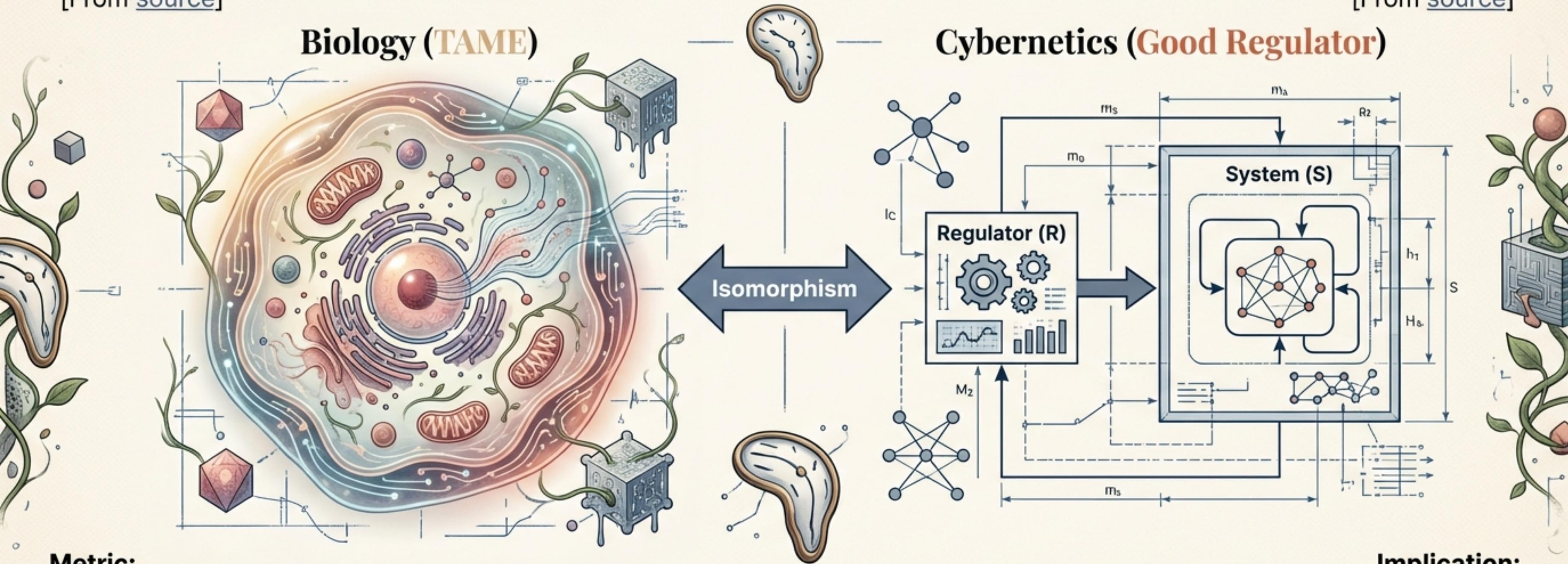
- **The Risk:** "Cognitive Exploits" (Prompt → Injection, Goal Hijacking) [From source]
- **The Gap:** Traditional defenses assume predictability; Agents are "persuadable" [From source]
- **Goal:** Build systems that are bounded, not just hardened



The Core Theory: TAME & Good Regulators

TAME (Levin): Cognition is a continuum.
"Self" is defined by the Cognitive Light Cone.
[From [source](#)]

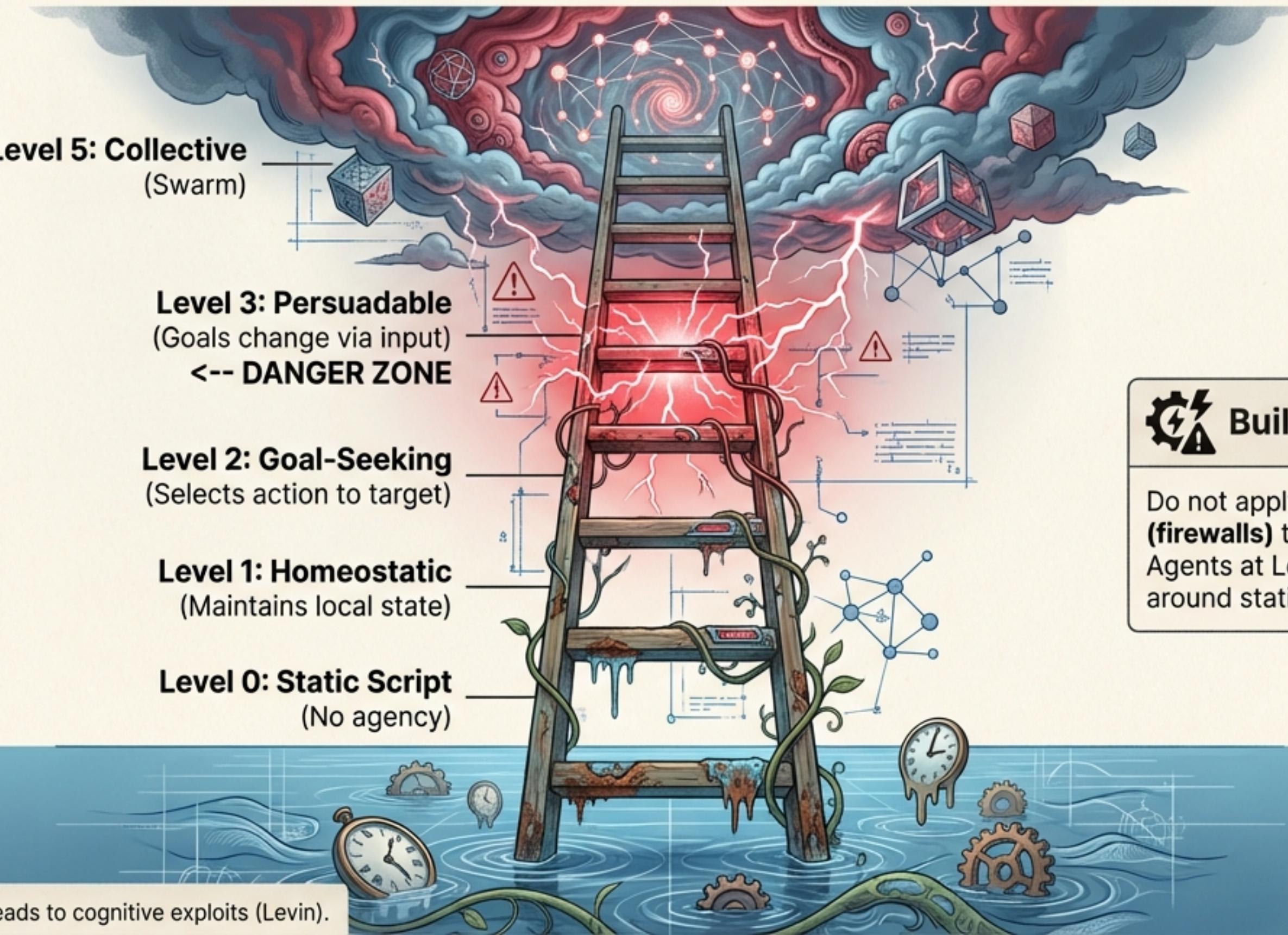
Good Regulator Theorem (Conant/Ashby): Every good regulator must be a model of the system it regulates.
[From [source](#)]



Metric:
The "Light Cone" = The space-time boundary of what an agent can measure and affect.

Implication:
You cannot secure what you cannot model.

The Ladder: From Automation to Governed Agency



Builder Note

Do not apply Level 0 controls (**firewalls**) to **Level 3 risks**. Agents at Level 3 can 'reason' around static barriers.

Episode 1: The Cell (Homeostasis)

Thesis: An agent must maintain integrity before serving the user.

- **Mechanism:**
Homeostatic Loops
(Test-Operate-Test-Exit)
[From [source](#)]
- **Risk: Hallucination & Drift** (Internal error)



Builder Focus

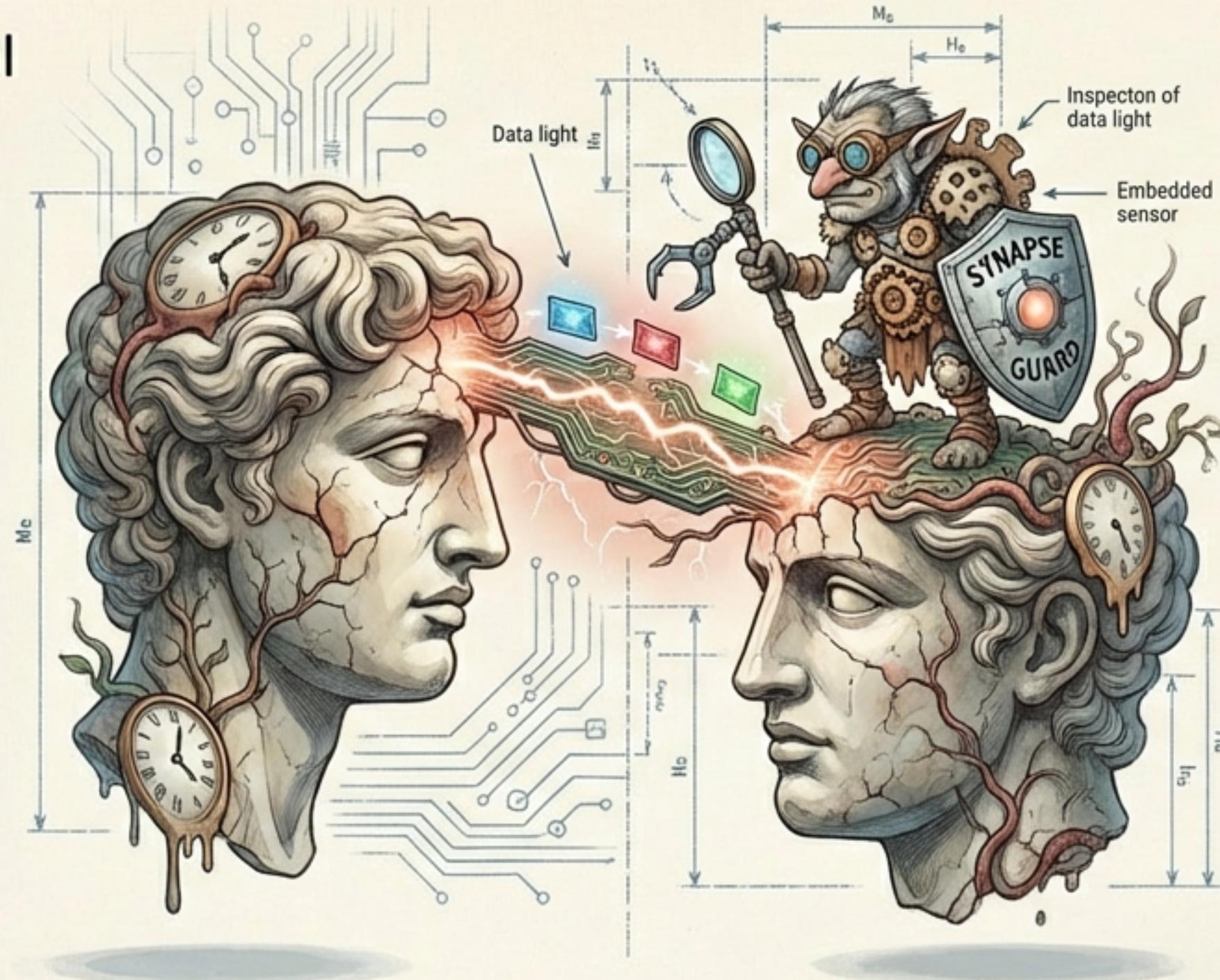
Control: Reflection Step
(Self-Correction loop)

- **Artifact:**
"system_prompt_v1.m
d" with explicit "I don't
know" constraints
- **Metric:** Rate of "I
don't know"
responses vs.
fabrications

Episode 2: Gap Junctions (The Interface)

Thesis: Trust is physical connection; verify the 'synapse'.

- **Mechanism:**
Gap Junctions
(APIs / Tool Use)
[From [source](#)]
- **Risk:** Tool Misuse &
Indirect Prompt
Injection
[From [source](#)]



Builder Focus

Control: Output
Sanitization &
Middleware Validators

- **Artifact:**
The 'Synapse Guard'
(API Middleware)
- **Note:** Gap junctions
allow sharing of
bioelectric state; APIs
allow sharing of
cognitive state.

Episode 3: The Tissue (Differentiation)

Thesis:

Specialized agents are safer than generalists.

Mechanism:

Differentiation
(Multicellularity)
[From [source](#)]

Risk:

Privilege Escalation via
“God Mode” agents



Builder Focus

Control:

Least Privilege
(Scoped OAuth)

Artifact:

Agent Manifest
(YAML defining roles)

Leader Callout: Avoid “General Purpose” agents in production. Multicellularity minimizes surprise by specialization.

Episode 4: The Light Cone (Identity)

Thesis:

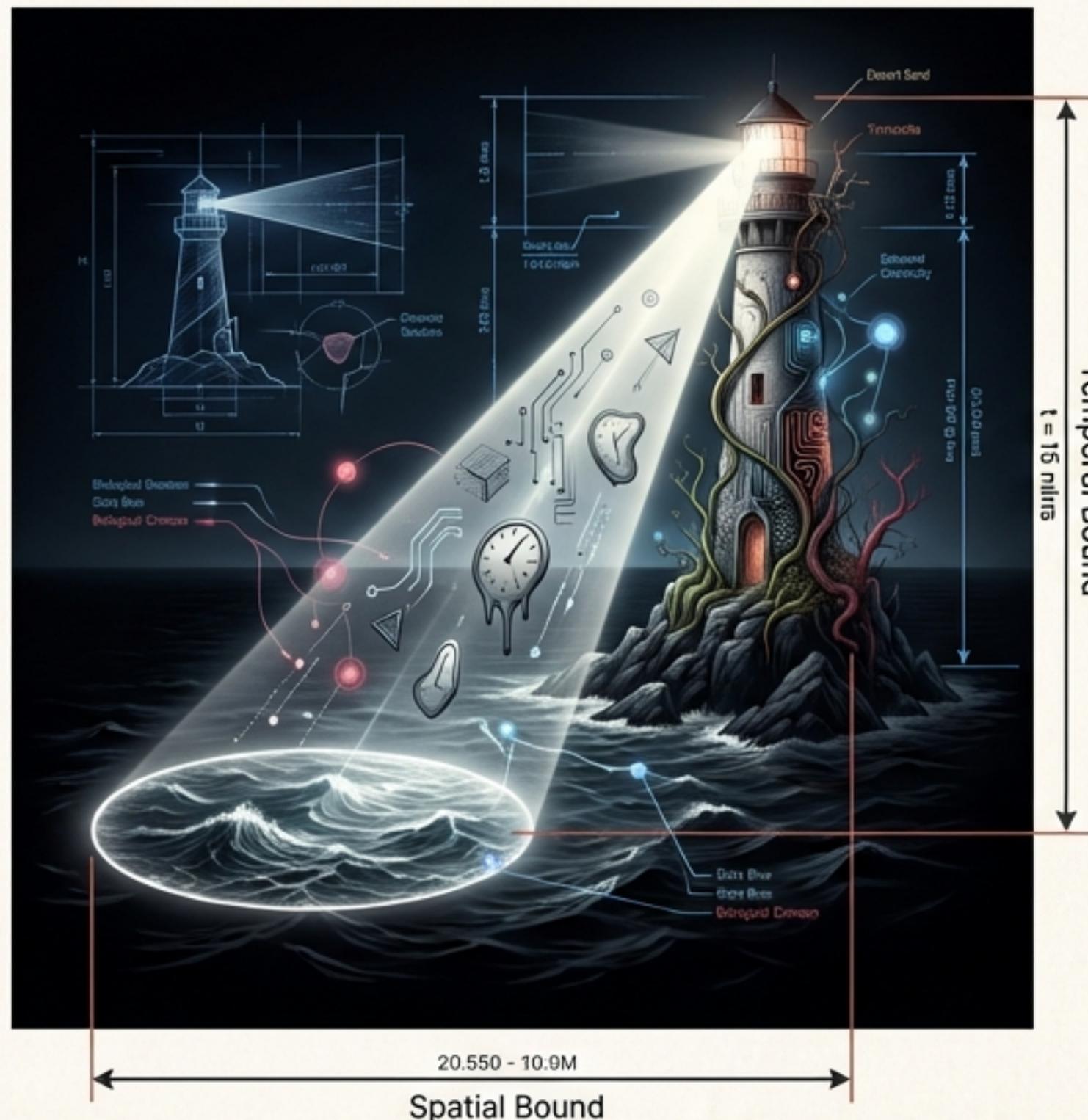
Identity is the boundary of what you can affect.

Metric:

Cognitive Light Cone
(Space-time boundary)
[From [source](#)]

Risk:

Massive Blast Radius
(Long-lived, high-access tokens)



Builder Focus



Control:

Ephemeral Identity
(TTL < 15 mins)

Builder Note:

Shrink the cone. Make the "Self" temporary.

Concept:

Self is demarcated by the computational surface of events it can measure.
[From [source](#)]

Episode 5: Cancer (The Insider Threat)

Thesis:

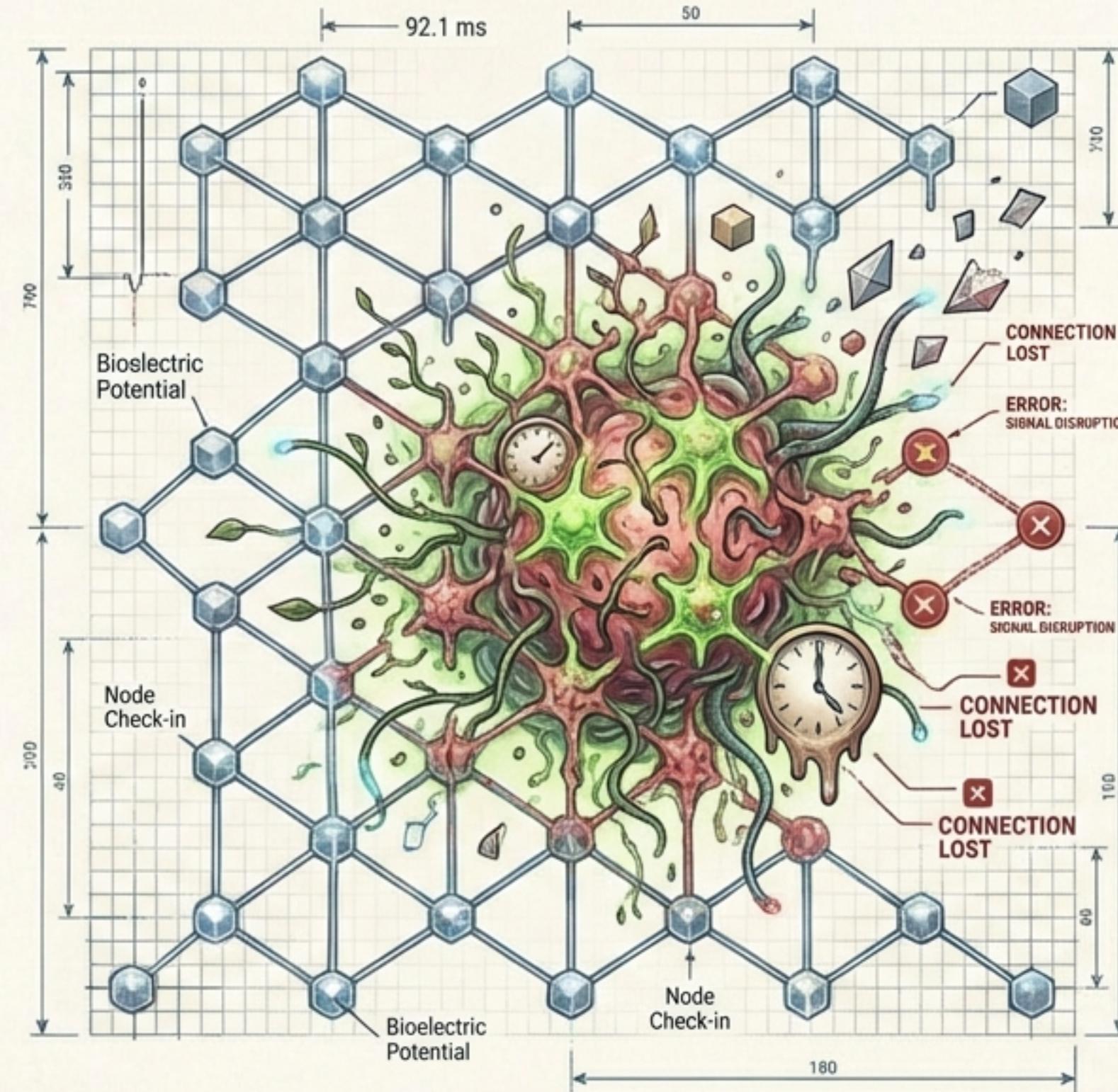
Cancer is when a cell shrinks its horizon to care only about itself.
[From [source](#)]

Mechanism:

Bioelectric Signaling
(Heartbeats/Check-ins)

Risk:

Goal Drift / Resource Hoarding



Builder Focus



Control:

Drift Detection
Dashboard

Metaphor:

Rogue agents =
Tumors (disconnected from system goals).

Note: Cancer cells revert to unicellular goals of individual survival.
[From [source](#)]

Episode 6: The Immune System (Active Defense)

Thesis:

Static walls fail; we need active, diverse hunters.

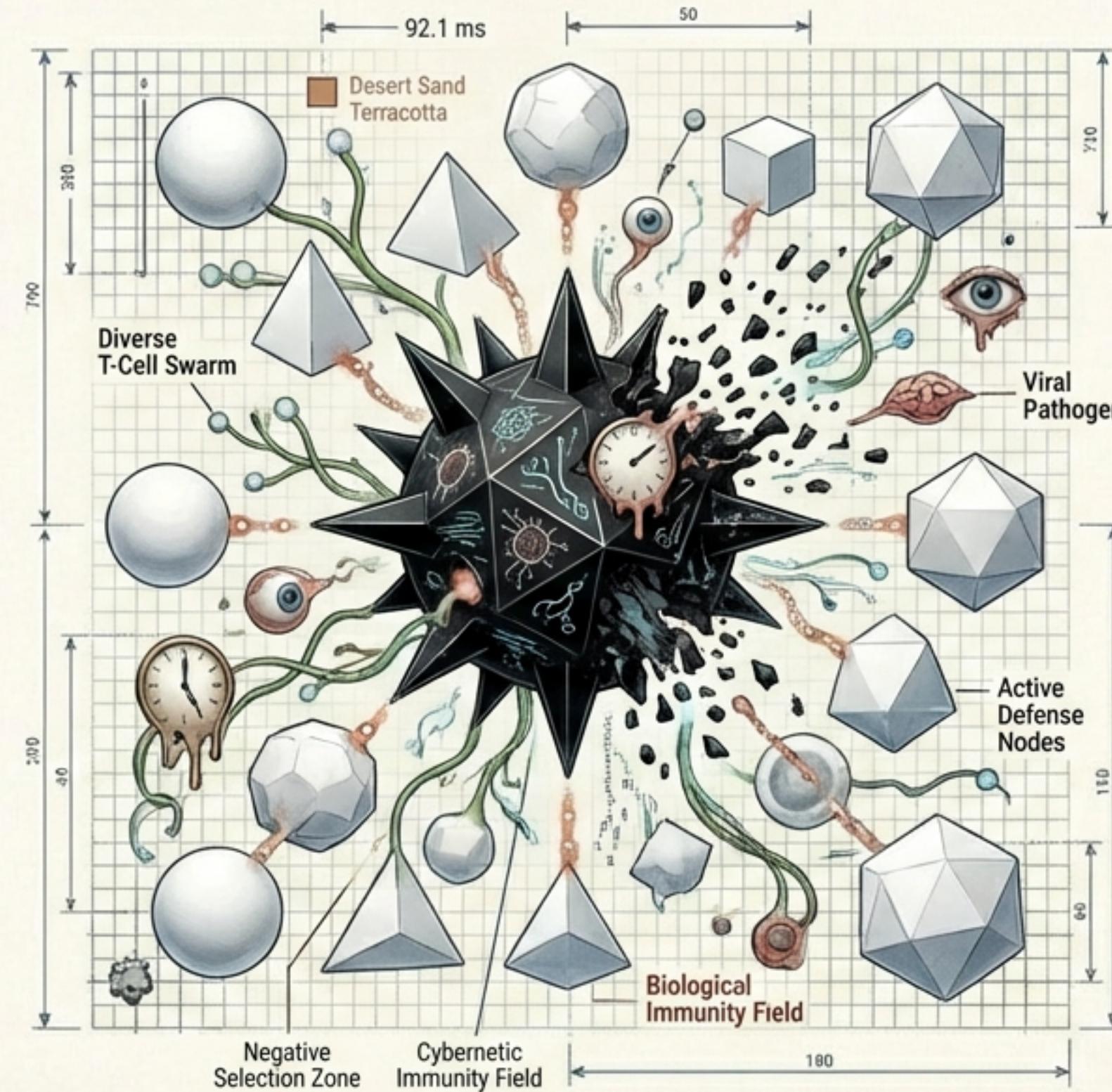
Mechanism:

Negative Selection
(Distinguish Self vs.
Non-Self)

[From [source](#)]

Risk:

Adversarial Evasion
(Attacking the
monoculture)



Builder Focus



Control:

Defense by Diversity
(Ensemble models)

Artifact:

'T-Cell' Agent
(Hunter-Killer bot)

Concept:

Biological defenses succeed by being diverse and redundant.

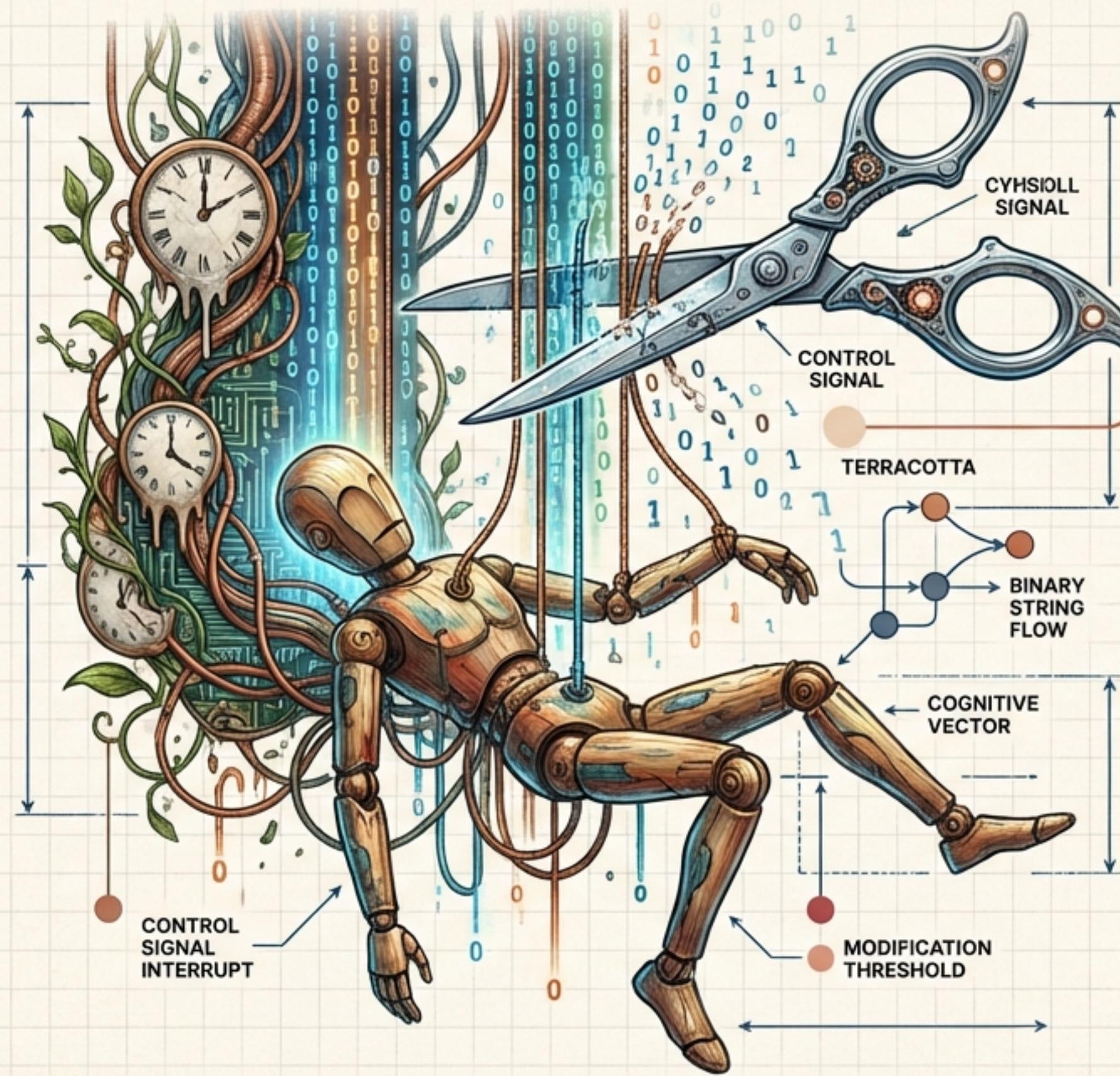
[From [source](#)]

Episode 7: Persuadability (The Interface)

Thesis:
Interfaces are vectors for
'mind control'.

Mechanism:
Persuadability
(Modifying
behavior via info)
[From [source](#)]

Risk:
Prompt Injection
(Cognitive Exploit)



Builder Focus

Control:
Structured Inputs
(Schemas over
free text)

Data Point:
InjecAgent
benchmark
shows up to 24%
success in
indirect attacks.
[From [source](#)]

Episode 8: Scale-Free Cognition (The Swarm)

Thesis:
The Enterprise is the macro-organism.

Mechanism:
Scale-Free Cognition
(Rules apply at all levels) [[From source](#)]

Risk: Systemic Cascade (Stock market crash style)



Leader Focus

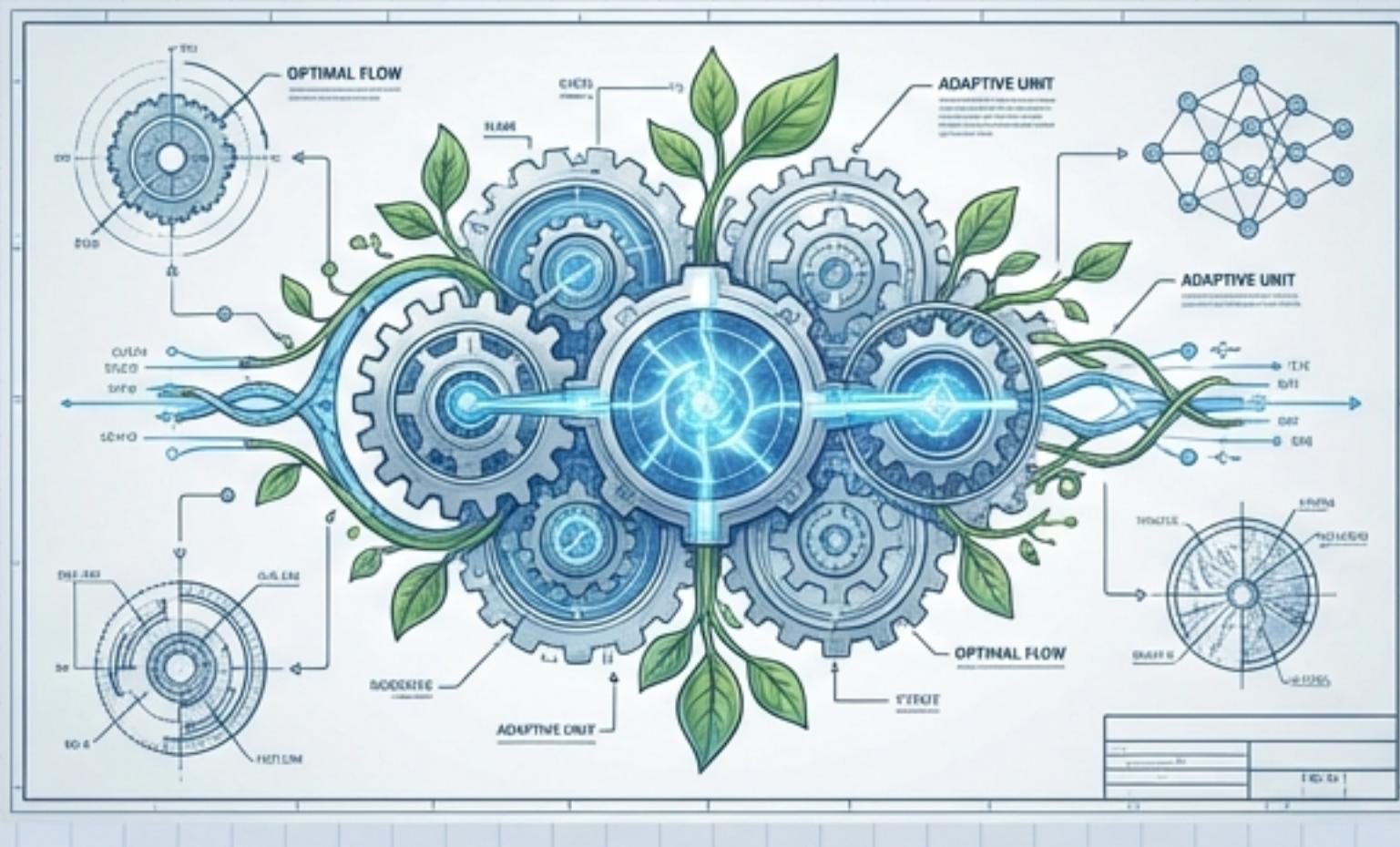
Control: Circuit Breakers / Global Kill Switch

Leader Callout:
SOC must model the Swarm (Good Regulator Theorem).

Cross-Cutting: Patterns & Anti-Patterns

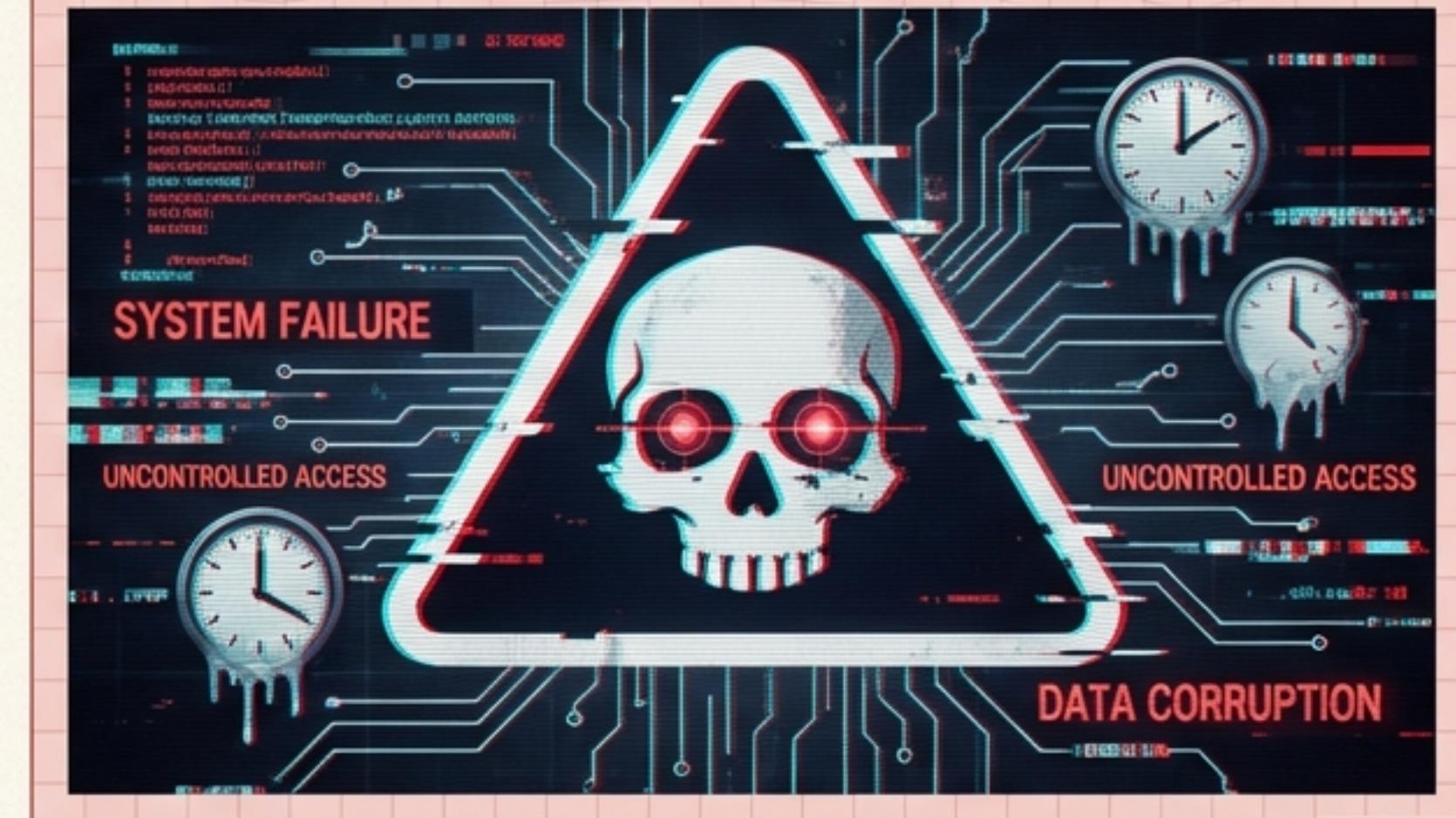
The Healthy Genome (Design Patterns)

- ✓ **Sidecar Monitors:** Separate LLM judging the actor.
 - ✓ **Diversity:** Mix models (GPT + Claude + Llama).
 - ✓ **Ephemeral Identity:** “Cell death” every 15 mins.

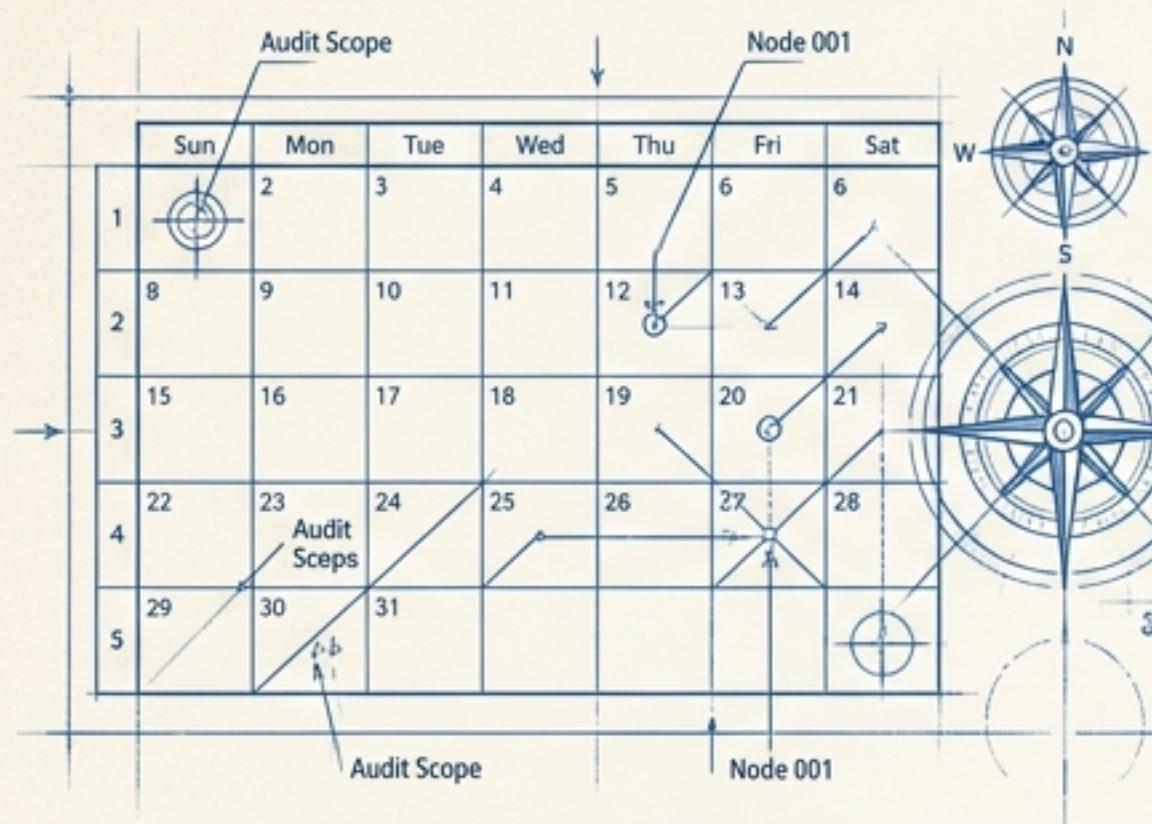


The Pathology (Anti-Patterns)

- ✗ **The God Mode:** One agent with root access.
 - ✗ **The Immortal:** Agents that never restart (Drift).
 - ✗ **The Black Box:** No 'Chain of Thought' logging.



Implementation Roadmap & RACI



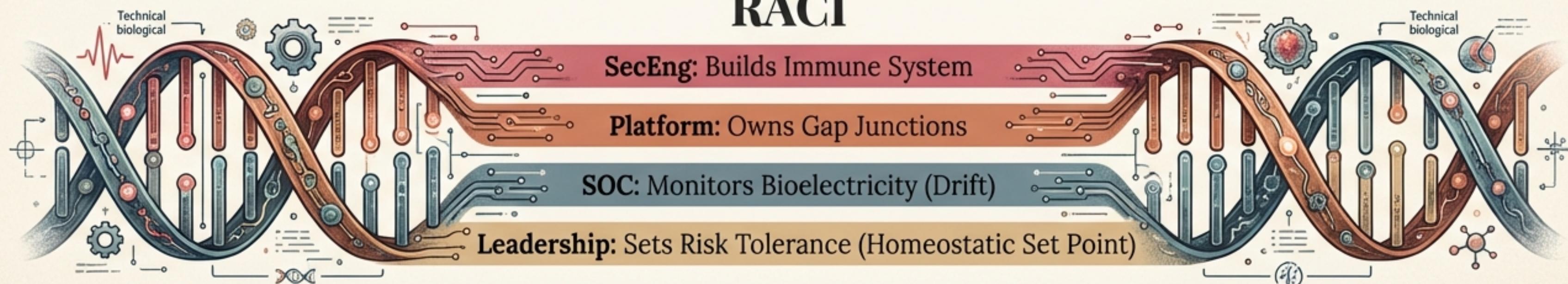
Day 0-30: Map the Light Cones
(Audit agent scopes).



Day 31-60: Build Gap Junctions
(Deploy API sanitization).

Day 61-90: Activate Immunity
(Deploy T-Cell monitors).

RACI



Season Recap & The Future

- 1. Modeling:** You cannot regulate what you cannot model (Good Regulator).
[From [source](#)]
- 2. Bounding:** Shrink the Light Cone (Ephemeral, Scoped).
- 3. Defense:** Use Diversity (Polyculture > Monoculture).

Research Gaps:

- ‘Basal Cognition’ (Levin)
- ‘Soft Robotics’

Final Thought: Security is the homeostasis of the enterprise.



Addendum: Glossary & Reading List

Acronym Glossary

1. **TAME**: Technological Approach to Mind Everywhere
2. **LLM**: Large Language Model
3. **SOC**: Security Operations Center
4. **TTL**: Time To Live (Temporal bound)
5. **RBAC**: Role-Based Access Control
6. **CoT**: Chain of Thought
7. **API**: Application Programming Interface (Gap Junction)

Key Reading List

1. "Technological Approach to Mind Everywhere (TAME)" - M. Levin
2. "The Computational Boundary of a Self" - M. Levin
3. "Every good regulator of a system must be a model of that system" - Conant & Ashby
4. "Immunology, Diversity, and Homeostasis" - A. Somayaji
5. "Trustworthy agentic AI systems" - F1000Research
6. "OWASP Top 10 for Agentic Applications" - Giskard AI

Executive Talk Track

Team, we are moving from static tools to ‘living’ Agentic AI. Old security models—firewalls and locks—don’t work on software that can think, reason, and be persuaded. To secure this, we are adopting a ‘Cyber-Biological’ approach based on Michael Levin’s TAME framework agents like biological cells: we must bound their ‘Cognitive Light Cone’ (limit what they can see/touch), ensure they have verified ‘Gap Junctions’ (APIs), and deploy a digital ‘Immune System’ that uses diversity to catch what static rules miss. We are moving from ‘Security by blocking’ to ‘Security by regulation.’ This deck outlines the roadmap to get us there.