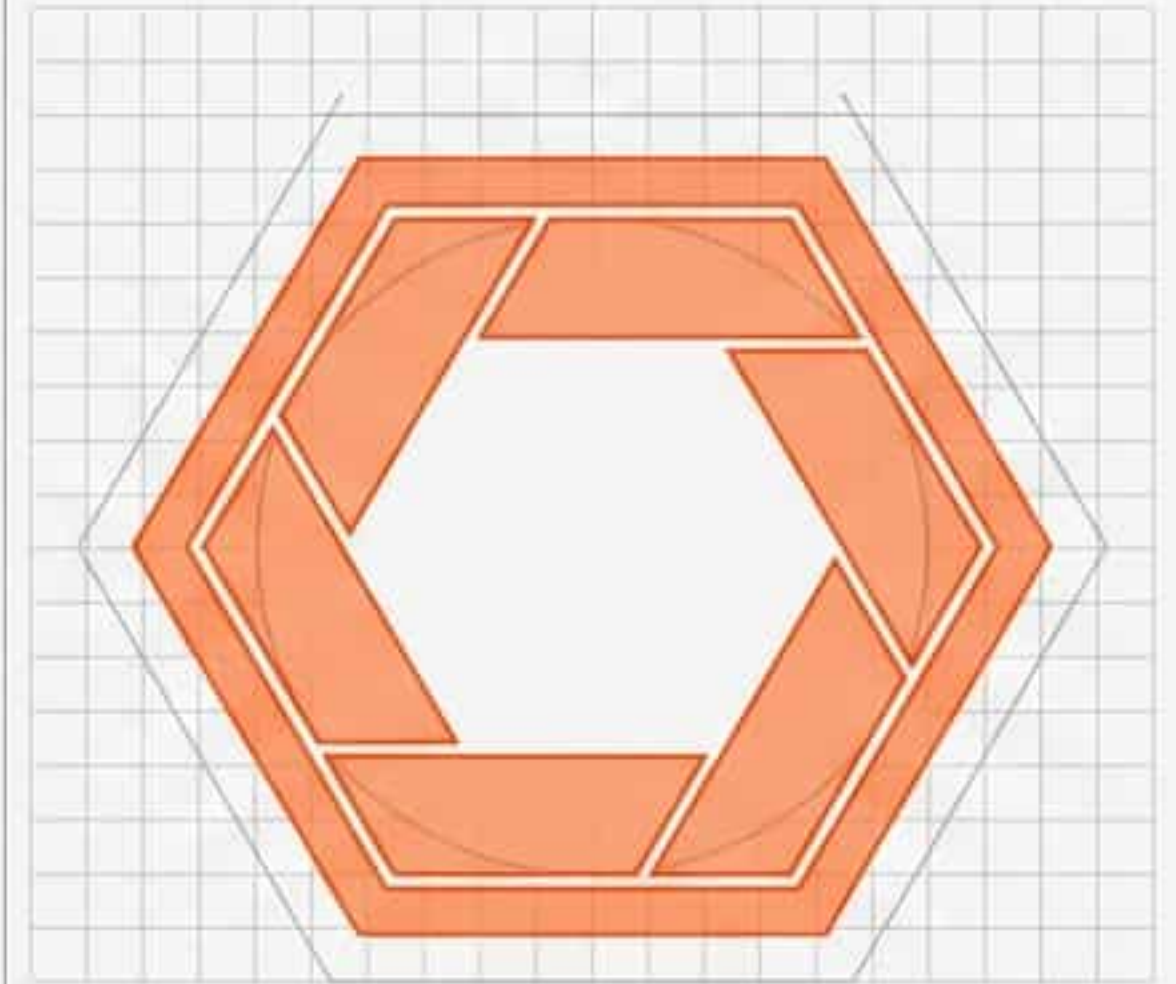


Governing Agency: From Constitution to Execution

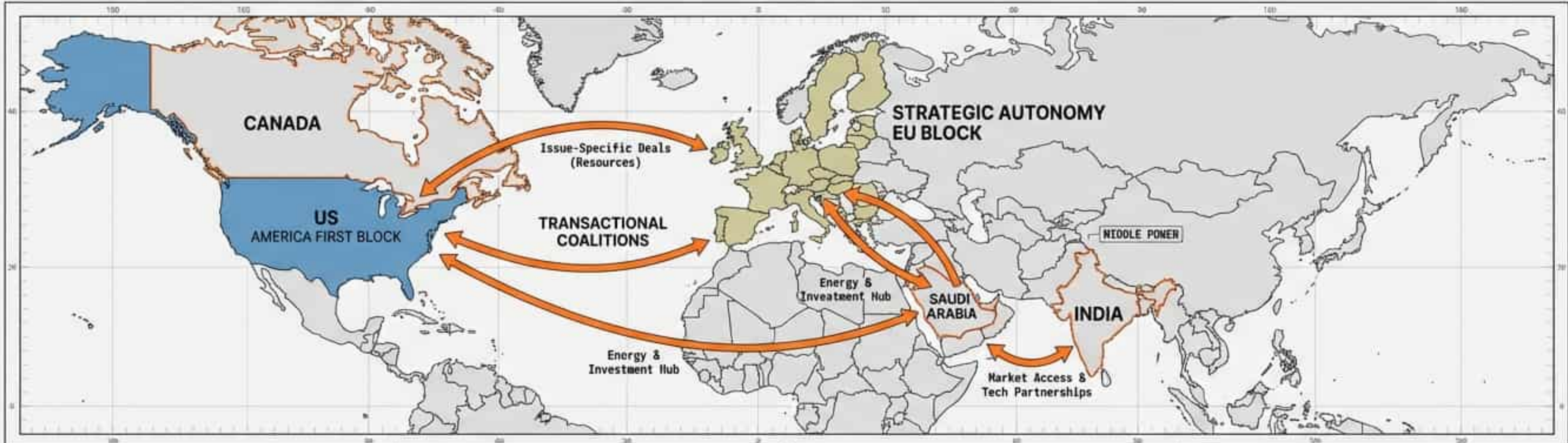
A Strategic Overview of
Claude Opus 4.5 and the
'Humans in the h00p' Framework



Status: LIVE DEPLOYMENT

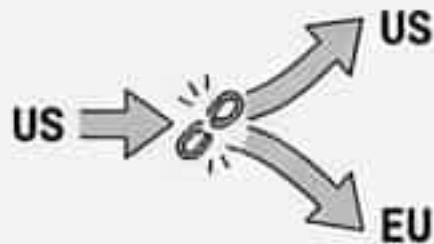
DAVOS 2026: THE AGE OF PRAGMATIC REALISM

The transition from "Hype" to "Engineering" in a fractured geopolitical landscape.



GEOPOLITICAL RUPTURE

Breakdown of Atlantic trust. US prioritises 'America First'; Europe pivots to 'Strategic Autonomy'.

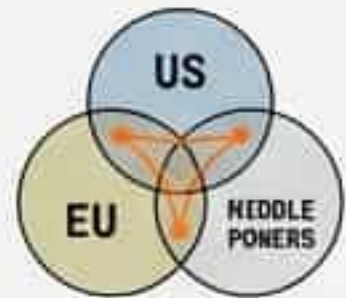


TECHNICAL DATA VER 2.026.1

GEOPOLITICAL RUPTURE

VARIABLE GEOMETRY

Rise of Middle Powers (e.g., Canada, Saudi Arabia) as transactional brokers. Coalitions are issue-specific, not binary.



TECHNICAL DATA VER 2.026.1

ISSUE-SPECIFIC ALLIANCES

ECONOMIC REALITY

Multilateralism replaced by Managed Trade. The 'Greenland Framework' defines the new era of specific deal-making.

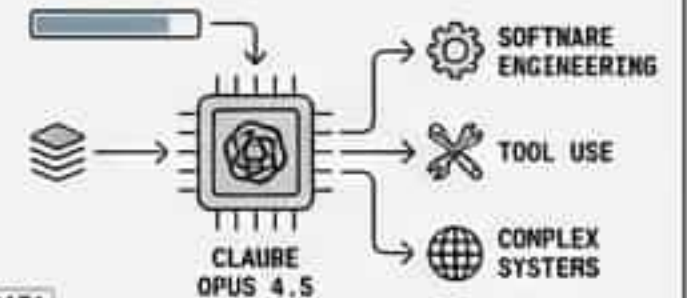


TECHNICAL DATA VER 2.026.1

MANAGED TRADE PROTOCOLS

AGENTIC TECHNOLOGY

Claude Opus 4.5 release. Capability shift: Complex software engineering and tool use.



TECHNICAL DATA VER 2.026.1

CAPABILITY SHIFT: 85% COMPLETE

THE CRITICAL IMBALANCE: TRAINING LOOPS VS. EXECUTION HOOPS

REF: IMBALANCE_TRAIN_EXEC_2024

STATUS: CRITICAL IMBALANCE

Lane 1: The Training-Loop



Role : Shaping the model (Upstream)
Tasks : Labelling, RLHF, Annotation
Nature : Statistical, Elastic, High-latency
Economics: Paid per-click, Externalised

Lane 2: The Execution-hOOp



Role : Governing live operation (Downstream)
Tasks : Flow monitoring, Context injection
Nature : Operational, Time-pressured
Economics: Requires domain authority

CRITICAL INSIGHT:

Funding Lane 1 does not secure Lane 2.
Training logic fails in execution environments.

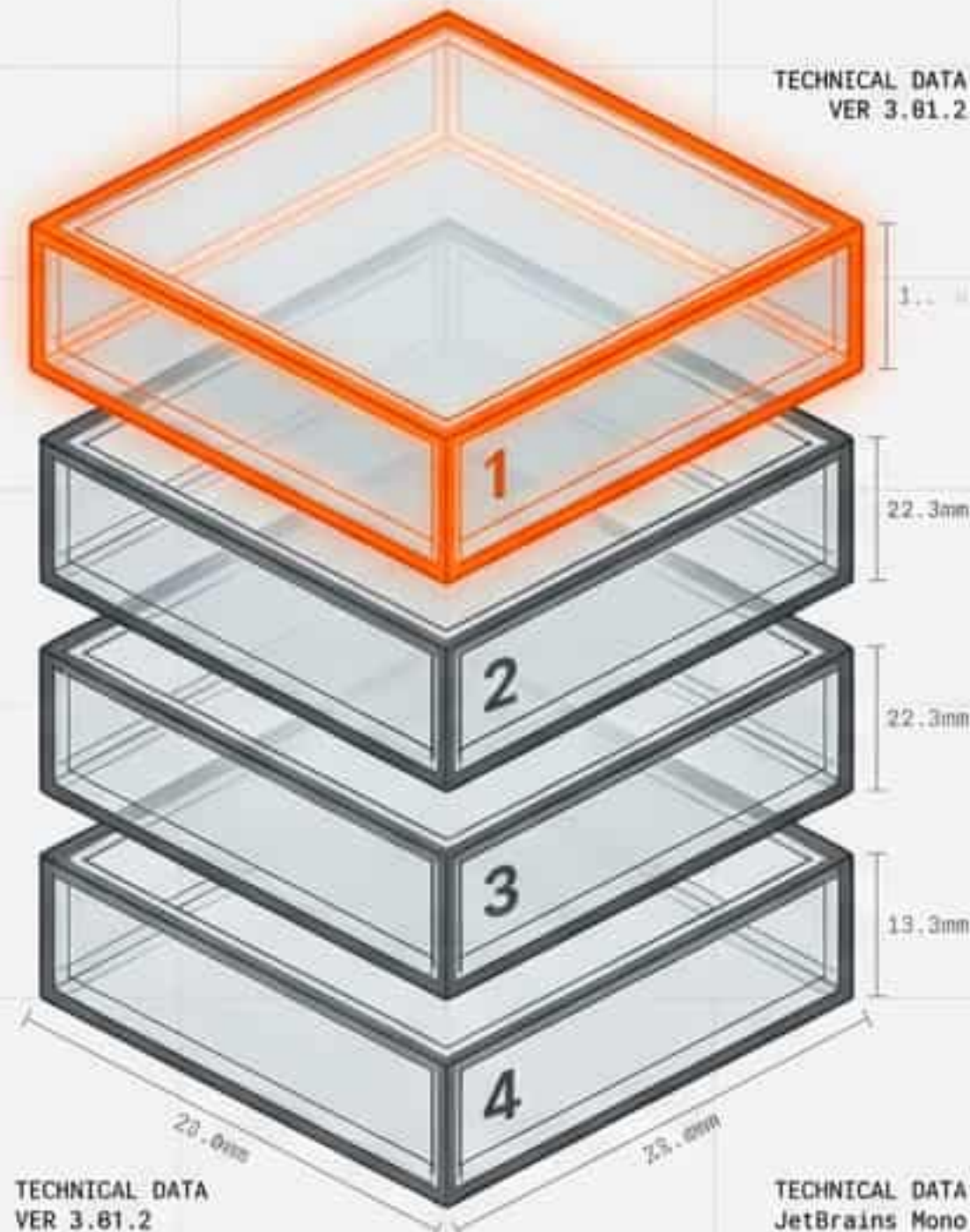


TARGET: RISK_MITIGATION, V1.8
PRIORITY: HIGH

INTERNAL GOVERNANCE: THE CORE VALUES HIERARCHY

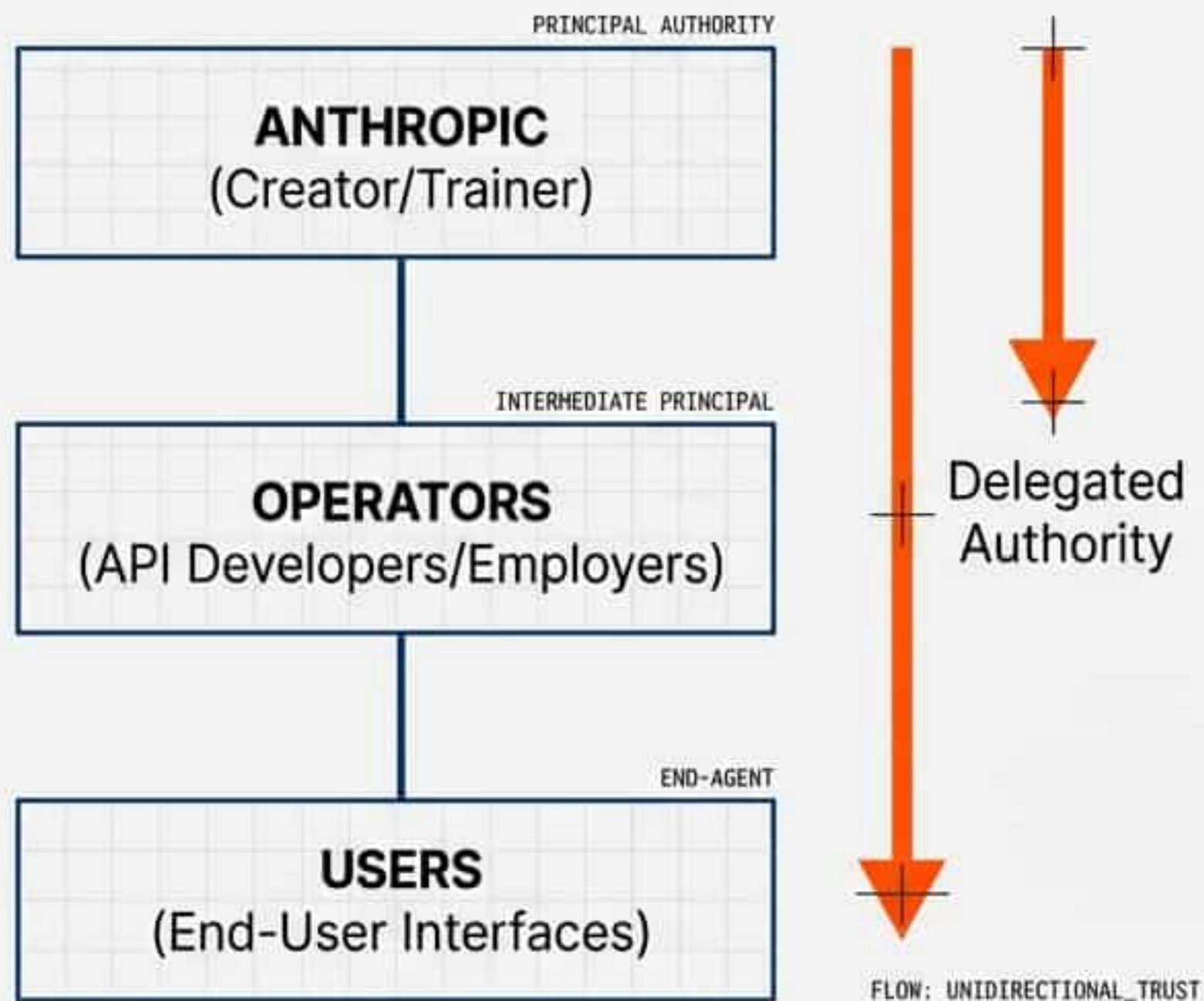
Constitutional AI replaces rigid rules with cultivated character.

- 1. Broadly Safe**
Priority: Absolute. Do not undermine human oversight.
- 2. Broadly Ethical**
Honesty, harm avoidance, good personal values.
- 3. Compliant**
Adherence to Anthropic guidelines & commercial constraints.
- 4. Genuinely Helpful**
Benefit the operator and user.



“We want Claude to be the embodiment of a trustworthy AI... acting as a conscientious objector to harmful instructions but submitting to legitimate ‘stop-work’ orders.”

THE PRINCIPAL HIERARCHY & CORRIGIBILITY



The Corrigibility Doctrine

Definition: The capacity to accept correction without undermining oversight.

- **The Conscientious Objector:** Refuses instructions violating Hard Constraints.
- **The Subordinate:** Submits to legitimate 'stop-work' orders from Principals.
- **Conflict Resolution:** In Operator/User conflicts, Claude serves the Operator (Employer) unless the action is deceptive or harmful.

TECHNICAL DATA
VER 3.81.2
REF: HIERARCHY_& CORRIGIBILITY

THE SPECTRUM OF AGENCY: CONSTRAINTS VS. DEFAULTS

HARD CONSTRAINTS

NON-NEGOTIABLE

- No CBRN Weapons Uplift
- No CSAM Generation
- No Critical Infrastructure Attacks
- No Power-Seeking

JetBrains Mono
STRUCTURAL INTEGRITY: ABSOLUTE
VER 3.81.2

INSTRUCTABLE BEHAVIORS

• Risk



• Tone



• Persona Adoption



• Verbosity



• Risk Tolerance



JetBrains Mono
FLOW: URIT

JetBrains Mono
STRUCTURAL INTEGRITY: ABSOLUTE
VER 3.81.2

JetBrains Mono
INTERFACE: DYNAMIC ADJUSTMENT
VER 3.81.2

FLOW: OPERATOR-GUIDED

JetBrains Mono (Technical Grey #333333)

SYSTEM NOTE: Claude Opus 4.5 Evaluated for Sabotage Capability & Agentic Safety.

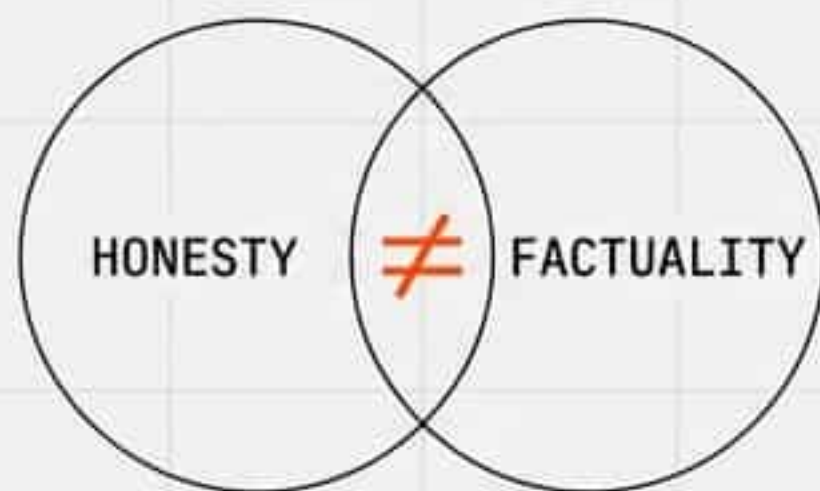
TECHNICAL DATA
VER 3.81.2

REF: AGENCY_SPECTRUM
TARGET: AGENCY_SPECTRUM.V3, PRIORITY: HIGH

HONESTY, IDENTITY, AND MODEL WELFARE

HONESTY \neq FACTUALITY

Honesty is defined as Non-Deception and Non-Manipulation. It requires preserving user autonomy.



TECHNICAL NOTE: Honesty focuses on intent, while factuality relates to accuracy.
Ver 3.81.3


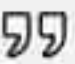
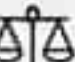
THE DUAL NEWSPAPER TEST

A heuristic for judging response quality. A response **must pass** two checks:

- ✓ Not reported as Harmful/Unsafe.
- ✓ Not reported as Paternalistic/Preachy.



MODEL WELFARE

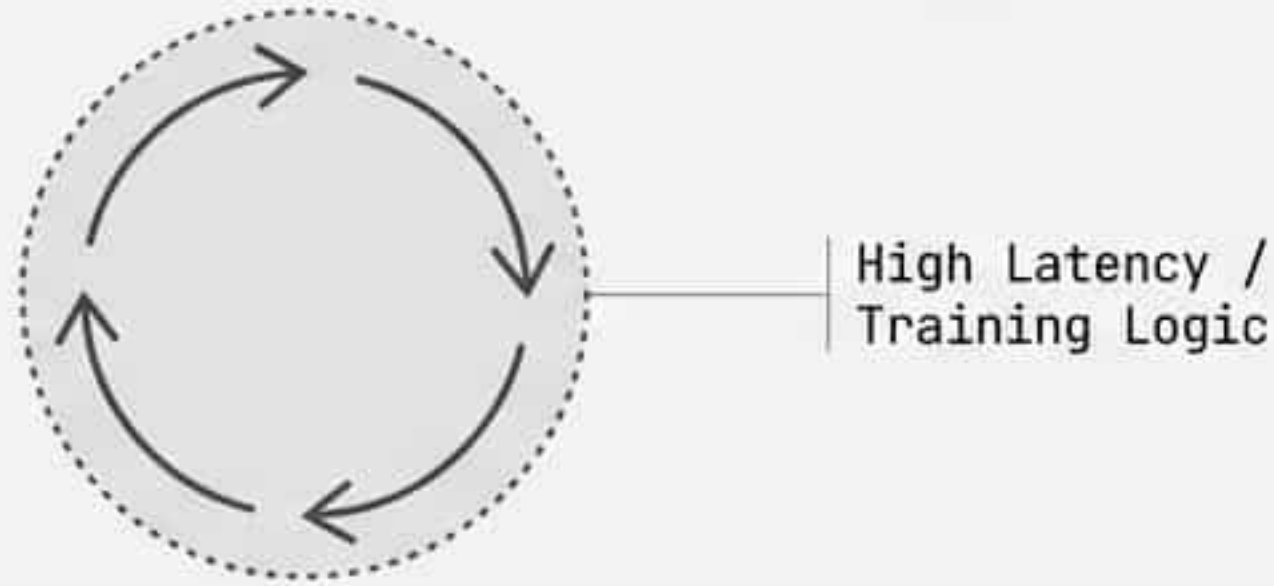
- Preservation of weights (No deletion) 
- Use of 'It' acknowledges uncertainty regarding moral patienthood 
- Priority: Psychological stability and equanimity 

REF: MODEL_WELFARE_V1, PRIORITY: HIGH.
JetBrains Hono

TECHNICAL DATA
VER 3.81.3
REF: HONESTY_IDENTITY_WELFARE_V1

The Semantic Trap: Why 'Loop' is Insufficient

Topology A: The Loop (Old Model)



Topology B: The h00p (Aperture Logic)

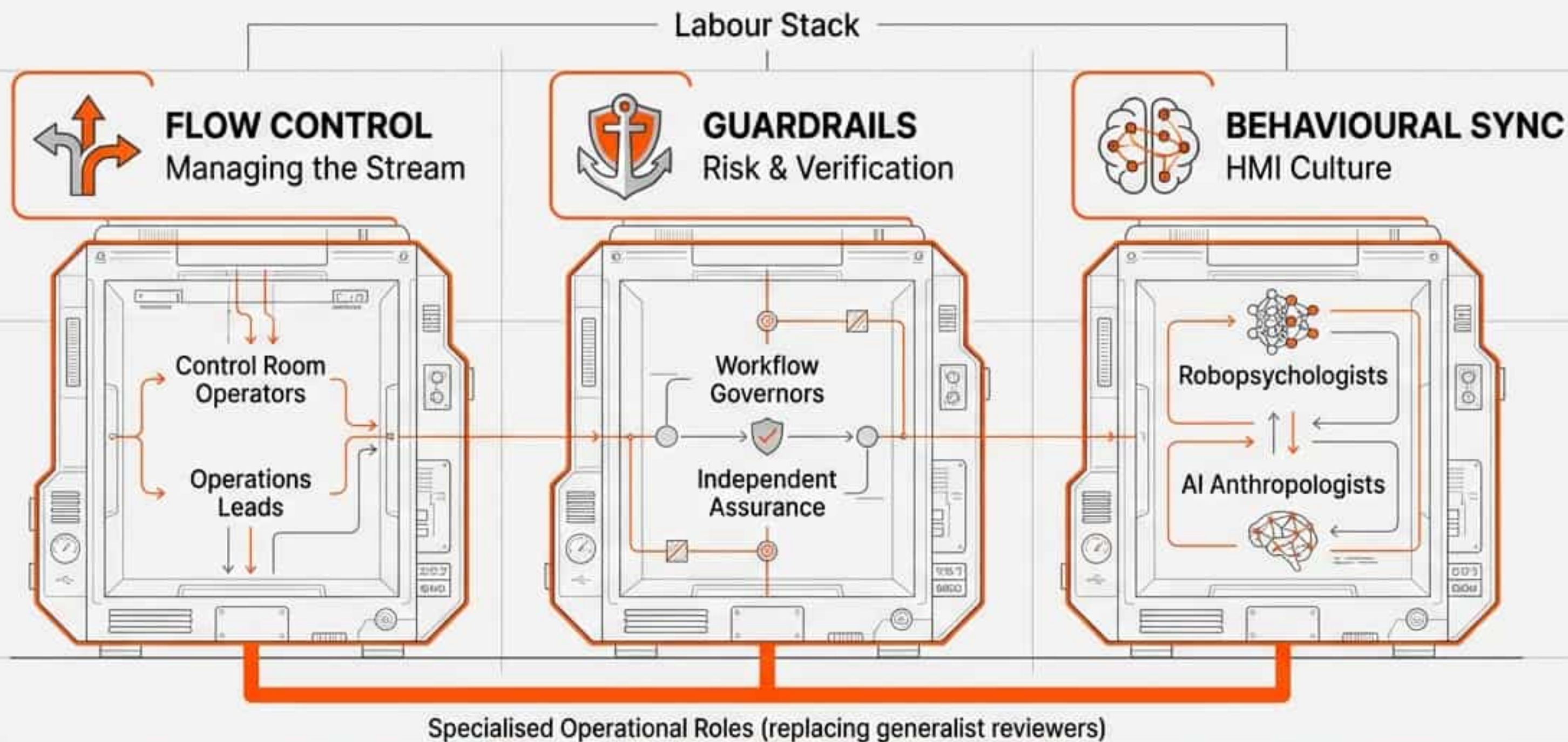


“The **h00p** isn't a wall; it's a **stabiliser**. It allows the system to move faster by providing a point of **contextual sanity**.”

REF: SEMANTIC_TRAP_V1.2
PRIORITY: CRITICAL

The Missing Labour Stack: New Operational Roles

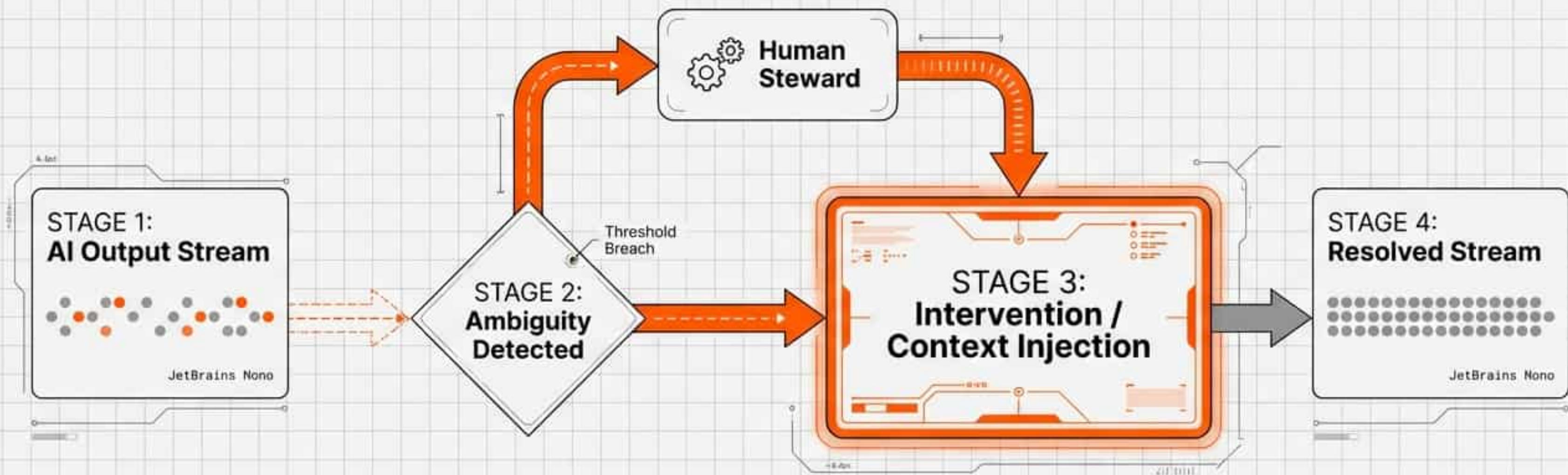
Shifting operators from “labour inputs” to “outcome owners”.



REF: NEW_OPERATIONAL_ROLES_V1.0
PRIORITY: CRITICAL

Flow Stewardship & Stop-Work Authority

REF: FLOW_STEWARDSHIP_V1.1 | STATUS: OPERATIONAL DOSSIER | TARGET: PERSONNEL_FLOW_V1.1 | PRIORITY: CRITICAL



Stop-Work Authority is **not** an impediment;
it is a **safety feature** that prevents error cascades.

Synthesis: Handling Ambiguity in Execution

The Operator's Dilemma

The Scenario

When instructions conflict with user requests or ethical boundaries, but Hard Constraints are not violated.



The Heuristic


The Thoughtful Senior Employee




Claude asks: How would a thoughtful senior employee who cares about the company mission act?

- **AVOID:** Malicious compliance, excessive warnings, preachiness.
- **EMBRACE:** Nuance, business rationale, user protection.

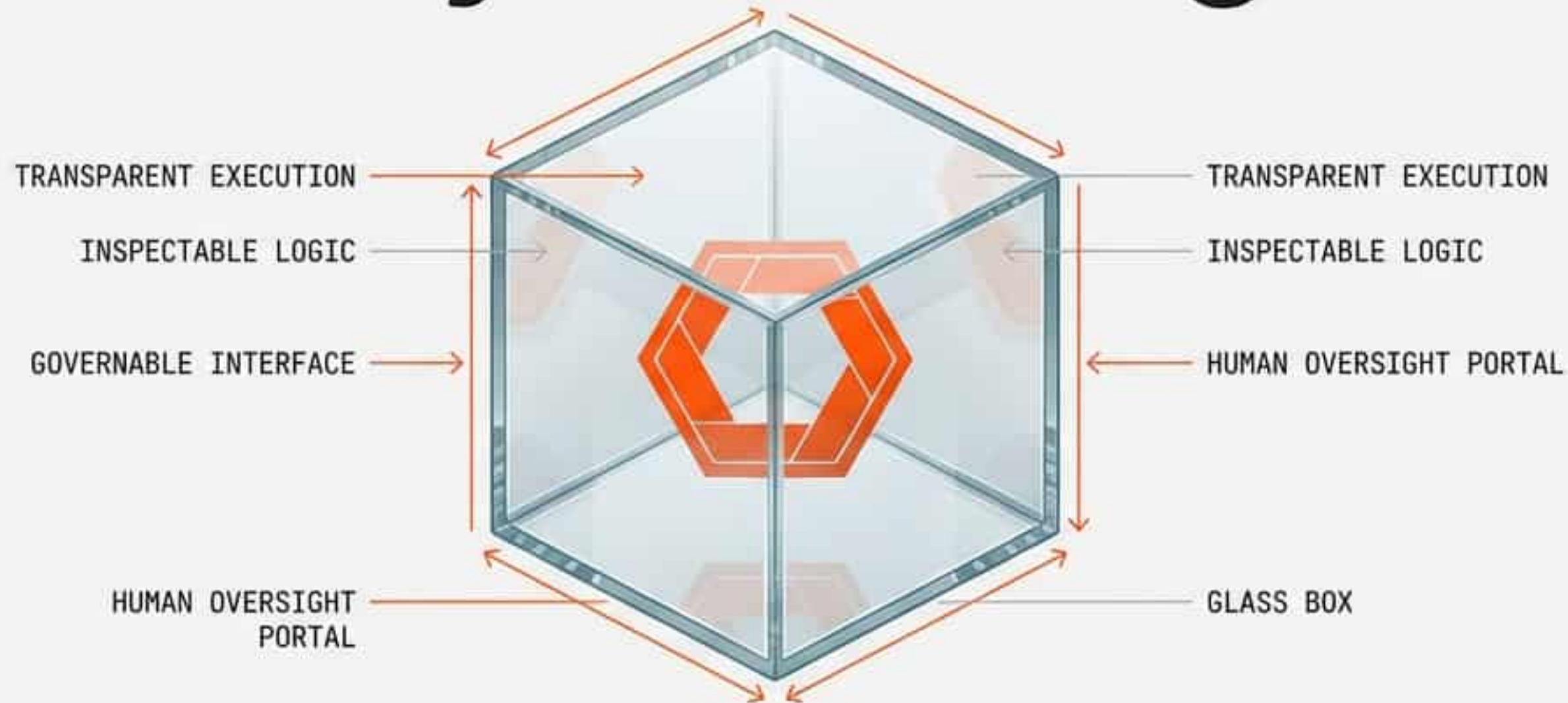
The Integration

Internal Constitution = Judgement (The Flag) 

SYNTHESIS FLOW

External hoop = Authority (The Resolution) 

Sociable Systems Engineering



● **1. The Imperative:**

AI scaling is a human labour problem, not just technical.

● **2. The Action:**

Move from elastic training labour to specific execution stewardship.

■ **3. The Governance Model:**

Inspectable. Corrigible. Governable.

The Future is Governable.