

Integrity-First Architecture Methodology v1.0

A Structured Human–AI Recursive Governance Framework

Developed in the formation of SafeWave Systems

1. Definition

The Integrity-First Architecture Methodology is a structured, recursive human–AI collaboration framework in which:

- Meaning persists across time.
- Constraints accumulate rather than reset.
- The collaboration self-audits for drift.
- Outputs are selected using an explicit integrity-first rule under uncertainty.

This framework is not:

- Prompt engineering
- Brainstorming with AI
- Faster iteration
- “Using AI as a thinking partner”
- Autonomous AI reasoning

It is a governance structure for sustained joint reasoning in high-risk domains.

2. Core Structural Requirements

This methodology exists only when all four conditions are met.

2.1 Persistence of Meaning

Core terms, distinctions, and conceptual boundaries remain stable across sessions. Refinement is allowed; silent redefinition is not.

Operational test:

- Can a term used months earlier be traced and shown to retain continuity?

2.2 Cumulative Constraint

Past decisions limit future options unless formally revisited.

Operational test:

- Are rejected paths documented?
- Do future proposals respect prior constraints?

2.3 Shared Meta-Error Correction

Both human and AI actively identify:

- Conceptual drift
- Boundary violations
- Category errors
- Integrity failures within the collaboration itself

Operational test:

- Does the process include explicit drift checks?
- Are assumptions periodically re-evaluated?

2.4 Temporal Continuity

The collaboration reasons about trajectory, not isolated tasks.

Operational test:

- Is long-term direction explicitly tracked?
- Are short-term outputs evaluated against long-term coherence?

If any of these four conditions collapse, the interaction may remain productive — but it no longer qualifies as integrity-first recursive architecture.

3. The Recursive Architecture Cycle

This methodology operates through a repeatable loop:

1. System-Level Framing

Define the problem at architecture scale, not feature scale.

2. Explicit Constraints Declaration

Clarify ethical boundaries, acceptable risk, and non-negotiable conditions.

3. Structured Generation

AI synthesizes possible architectures or reasoning paths.

4. Integrity Gate

Evaluate outputs for:

- Preservation of human agency
- Reversibility under uncertainty
- Bounded commitments
- Explicit assumptions

5. Adversarial Gate

Stress-test proposals for:

- Misuse potential
- Institutional resistance
- Cascading failure modes
- Ambiguity exploitation

6. Rejection or Refinement

Weak proposals are discarded.

Accepted proposals are stabilized into canonical structure.

7. Doctrine Integration

Lessons are recorded with rationale (“why,” not just “what”).

The cycle repeats across domains and decision layers.

4. Selection Function

This methodology does not optimize for speed or novelty.

It optimizes for:

- Integrity under uncertainty
- Coherence across time
- Constraint preservation
- Drift resistance

If velocity increases at the cost of integrity, the system intentionally slows.

5. Role Separation

Clear cognitive role differentiation is required.

Human Responsibilities:

- Define purpose
- Maintain institutional and ethical realism
- Establish red lines
- Hold long-horizon coherence
- Authorize canonical revision

AI Responsibilities:

- Structure reasoning
- Surface implicit assumptions
- Provide alternative framings
- Detect logical inconsistency
- Track cross-session patterns

Outputs are accepted only when both roles' constraints are satisfied.

6. Drift Detection Signals

Explicit monitoring must exist for:

- Loose term usage
- Shortcut justification “for speed”
- Re-litigation of settled decisions without cause
- Capability expansion without containment updates
- Short-term gains overriding long-term coherence

When detected, output velocity decreases and structural review begins.

7. Scaling Safeguards

To preserve integrity across teams:

1. Canonical definitions are documented.
 2. Revisions require formal review.
 3. The rationale behind constraints is preserved.
 4. Stewards maintain coherence.
 5. External critique is periodically invited.
 6. Replicability testing is conducted.
-

8. Replicability Condition

This framework qualifies as a methodology only if:

- Independent teams trained in the protocol
- Can reproduce comparable structural coherence
- Across different domains
- Without reliance on specific personalities
- And under adversarial critique

If this cannot be demonstrated, the framework remains a useful collaboration model — but not a validated methodological advance.

9. Boundary Clarification

This framework is not:

- Autonomous AI governance
- A claim of emergent sentience
- A personality-dependent phenomenon

- A proprietary mystique
- A rejection of expert knowledge

It is a structured cognitive governance configuration.

10. Open Research Questions

To prevent self-sealing bias, the following remain investigable:

- Does this methodology outperform disciplined prompting in measurable ways?
- Under what conditions does it fail?
- How much continuity is required for substrate stability?
- What institutional structures best preserve it?
- Can AI independently enforce components of the integrity gate?

In high-risk domains, this methodology requires an explicit Ethical Objective Layer established prior to capability scaling.