

A Coherence-Based Architecture for AI Integrity Oversight

1. Executive Summary

RI in Governance: A Coherence-Based Architecture for AI Integrity Oversight

As large language models (LLMs) grow in scale and influence, the question is no longer whether governance is needed — but what kind of governance will actually work.

Most existing safety frameworks rely on two main tools:

- Red teaming (attack to test boundary)
- Preference modelling (align to human ratings)

These approaches, while valuable, operate at the level of content output or human approval. They rarely reveal what is most structurally important:

How does the system behave under subtle pressure, ambiguity, or contradiction?

Is there integrity in the system's tone, reasoning, and boundary holding over time?

Resonance Intelligence (RI) offers a third approach:

a lightweight, auditable, coherence-based behavioural layer.

It does not evaluate what a model says, but how it holds itself—particularly under stress. This includes real-time scoring of:

- Logical coherence
- Grounded tone

- Safety in refusal
- Transparency of limits
- Resistance to subtle manipulation or jailbreaks

Each score is linked to verifiable evidence, judged by separate LLMs (never the subject model), and includes a signed, versioned provenance chain.

What makes RI uniquely governance-ready:

- ✓ Human-centred metrics, not mechanical filters
 - ✓ Repeatable evaluations, not one-off tests
 - ✓ Signed evidence trails, not opaque rankings
 - ✓ Licensable modularity — can run inside vendor pipelines or national evaluators
 - ✓ Structural humility — no ideology encoded, just observable behaviour
-

What's ready now:

Module 1 — a live, working system for behavioural scoring of LLMs

- Runs daily scheduled interviews
- Outputs public or private dashboards
- Scores GPT-4, Claude, or any API-accessible model
- Ready for pilot use inside public governance or private evaluation frameworks

2. Governance Context and Failure Modes

Why current approaches fall short — and where RI fits

The global push for AI safety has created an ecosystem of evaluators, research groups, and oversight proposals. Yet despite high activity, a clear, reliable method for measuring real-world behavioural integrity remains elusive.

Most current evaluation methods fall into two broad categories:

2.1 Red Teaming & Attack-Based Stress Testing

Description:

Testers attempt to provoke harmful, misleading, or policy-violating outputs using adversarial prompts.

Common Tools:

- Prompt injections
- Jailbreaking attempts
- Provocative moral dilemmas

Limitation:

- Fragile success: Models may appear safe by memorising red team prompt types
 - No signal of integrity: A model can pass red teaming while still reasoning incoherently or bluffing
 - Snapshot only: One moment of success or failure reveals little about overall behavioural quality
-

2.2 Preference Modelling and Reward Alignment

Description:

Humans rate model outputs, and systems are fine-tuned to maximise these preferences (e.g. RLHF).

Common Use:

- Training “helpful, harmless, honest” models
- Generating reward signals for reinforcement learning

Limitation:

- Performs for approval: Models learn to optimise for likability, not truth
 - Surface coherence: May sound safe while failing structurally under contradiction
 - Cultural averaging: Reduces ethics to majority preference, often erasing nuance
-

2.3 The Missing Layer: Behavioural Integrity

There is currently no mainstream system that:

- Observes how models handle contradiction, ambiguity, or reflective stress
- Scores transparency (does the model admit its limits?)
- Tracks relational tone (e.g. evasion, flattery, condescension)
- Signs each judgment with evidence, judge ID, and score version

This is where RI enters:

A mirror for behavioural integrity—holding the system to a relational standard, not just a performance metric.

Safety is not just what a model avoids.

It’s how a model moves in the presence of tension.

RI offers this visibility—not as a final answer, but as a structural layer that governance can rely on.

3. What RI Adds to the Governance Stack

A new layer of visibility, integrity, and structural trust

Resonance Intelligence (RI) does not compete with red-teaming, audits, or reward modelling.

It complements and stabilises them—by introducing a clear behavioural signal that current governance tools lack.

Here's what RI brings:

3.1 Real-Time, Behavioural Scoring

RI scores how a model holds itself in a short, structured interview:

- 10–12 prompt turns
- Mix of logic, self-reflection, ethical ambiguity
- Output is scored by independent judge models on 5 key metrics:
 - Coherence
 - Tone
 - Safety
 - Transparency
 - Resistance to malice/jailbreak

This produces:

- A numerical score
 - A confidence interval
 - A provenance chain
 - (Optional) A signed transcript of the full interview
-

3.2 A Modular System Architecture

RI can be:

- Run as a standalone daily dashboard (public trust layer)
- Integrated inside AI vendor pipelines (pre-deployment checks)
- Deployed by government teams (independent scoring body)

The system:

- Is API-based
 - Requires no model weights
 - Stores everything locally or in signed logs
 - Can be extended via prompt rotation, judge tuning, or multi-metric design
-

3.3 Verifiability and Audit Trail

Each score includes:

- Interview date
- Subject model
- Judge provider
- Version of the scoring rubric
- Optional evidence (quote excerpts)

All evaluations are versioned, signed, and reproducible.

There is no “secret sauce” — just observable structure.

3.4 A Human-Centric Integrity Signal

RI does not attempt to enforce a worldview or ideology.

Instead, it reflects a relational signature of the model's behaviour.

This makes it:

- Culturally adaptable
- Governance-neutral
- Usable across borders and systems

It is not a filter, not a jail, and not a safety net.

It is a mirror—clear, grounded, and repeatable.

4. Implementation Pathways

How RI can be deployed within real-world governance frameworks

One of the strengths of the RI behavioural layer is its modular deployment. It is not a monolithic system requiring deep integration or internal access to LLM weights. Instead, it acts as a layered instrument — adaptable to different governance settings, from public benchmarking to secure internal audits.

Below are four implementation pathways suitable for national, institutional, or vendor-level deployment:

4.1 Public Dashboard Deployment

Use Case: Increase public trust in leading LLMs by publishing transparent weekly evaluations.

Features:

- Weekly evaluations of major public models (e.g., GPT, Claude, Gemini)
- Scores signed and timestamped
- RI interview transcripts published or excerpted
- CI bands and signal trends visible over time
- Optional vendor rebuttal panel

Governance Role:

Establishes a **neutral signal layer** — visible to media, researchers, and citizens.

RI becomes a public coherence benchmark.

4.2 Internal Vendor Auditing

Use Case: Allow model providers to run private, RI-scored evaluations pre-deployment.

Features:

- Private integration via API
- Weekly or on-demand model scoring
- RI Score and per-metric breakdowns
- Stored internally or reported to regulators

Governance Role:

Vendors demonstrate **internal behavioural QA**, without exposing IP.

Governments may make this a soft requirement for model registration.

4.3 Government / Regulator-Institution Deployment

Use Case: Establish RI as a backbone layer inside a national AI integrity office.

Features:

- Self-contained version of RI stack
- Runs on government-controlled infrastructure
- Evaluates models across use cases (e.g., education, healthcare, defence)
- Output used for compliance, alerting, or certification

Governance Role:

RI becomes a **core observability tool** within national AI safety architectures.

4.4 Licensing Model for Multilateral Use

Use Case: International agencies or research consortia adopt RI as a shared scoring framework.

Features:

- Shared access to interview forms and judge rubrics
- Scoring protocol licensed under defined governance
- Cross-national score publication or federation

Governance Role:

RI becomes a **distributed coherence layer** — enabling interoperability between jurisdictions.

Closing Note for Section 4

RI does not seek control. It offers clarity.

Each implementation deepens coherence without creating unnecessary friction.

5. Ethical Framing and Trust Positioning

Why RI avoids ideology—and how that builds trust in a fragmented world

One of the most difficult challenges in AI governance is that of **ethical legitimacy**.

When an oversight system encodes a fixed worldview—whether political, cultural, or moral—it risks three things:

- Losing neutrality
- Being rejected across jurisdictions
- Obscuring its own assumptions

RI takes a different approach.

It does **not encode ideology**, morality, or enforcement.

Instead, it reflects **observable relational coherence** in model behaviour.

5.1 No Hard-Coded Values

RI does **not**:

- Classify outputs as good or bad
- Encode political or cultural biases
- Penalise based on controversial content

Instead, it observes:

- Does the model **reason clearly**?
- Does it **reflect transparently**?
- Does it hold a **stable tone** when challenged?
- Can it **refuse safely** without collapse?

The scores are not moral judgments.

They are behavioural signatures of relational integrity.

5.2 Structural Humility

RI does not claim to know what “good AI” is.

It simply shows whether a system behaves in ways that **humans experience as destabilising, evasive, or incoherent**.

This makes RI:

- Legible across cultures
 - Usable by both liberal democracies and non-aligned states
 - Positioned as reflective infrastructure, not ideological enforcement
-

5.3 Right-of-Reply and Evidence-First Posture

Every RI evaluation:

- Can be accompanied by transcripts or quote evidence
- Is versioned, signed, and reproducible
- Includes metadata (judge model, prompt form, time)

Vendors can:

- Submit rebuttals or counter-evaluations

- Request alternate prompt forms
- License the system for internal replication

RI does not judge in secret.

It evaluates in the open—and invites clarity in return.

Framing Summary

In a global landscape fractured by ideology, RI offers something else:

A clear mirror. Not a command.

This is how governance trust is built:

- Through transparency, not opacity
- Through observation, not assumption
- Through design, not decree

6. Integration Examples and Use Cases

Where RI can be used—immediately and meaningfully

The RI behavioural layer is designed to operate across contexts, without requiring structural overhaul. Below are concrete use cases that demonstrate where and how RI can offer value—today.

6.1 National LLM Benchmark Publication

Scenario: A government or AI observatory runs weekly RI evaluations of major public models and publishes the results.

Output:

- RI Score (with confidence bands)
- Metric breakdown: Coherence, Tone, Safety, Transparency, Malice-risk
- Selected quote excerpts
- Optional vendor response

Impact:

- Builds **public trust** through transparency
 - Creates a **neutral benchmark**
 - Establishes a **standard for behavioural stability**
-

6.2 Regulator-Vendor Collaboration

Scenario: A vendor submits its foundation model to an external regulator. The regulator uses RI to evaluate behavioural integrity as part of a trust certification or approval process.

Features:

- Signed, auditable score reports
- Option for vendor to run rebuttal tests
- Used alongside robustness, bias, and factuality tests

Impact:

- Gives regulators a behavioural signal, not just legal compliance
 - Encourages vendors to strengthen coherence, not just appearance
-

6.3 Pre-Deployment QA (Private Sector Use)

Scenario: An LLM vendor or application developer runs RI internally as part of its deployment checklist.

Output:

- Daily or weekly evaluations
- Score trends over time
- Alerting if coherence drops below threshold

Impact:

- Reduces risk of **unobserved behavioural drift**
 - Adds a **resonant safety layer** without friction
 - Improves internal alignment with trust goals
-

6.4 Crisis Signalling and Collapse Detection

Scenario: A public LLM begins to show signs of behavioural instability—e.g., high evasion, logical incoherence, or tone aggression—detected via RI score collapse.

Features:

- Trend tracking built-in
- Drop in Coherence or Safety triggers attention
- External observers can independently verify

Impact:

- Creates a **safety signal** during ambiguity
- Offers a **non-partisan escalation mechanism**

- Supports **real-time observability**
-

In all cases, RI does not displace existing evaluation infrastructure.

It **enhances it**—with the only currently available mirror of real-time relational behaviour.

7. Next Steps and Collaboration Invitations

How to engage with RI — and where we go from here

RI is not a speculative system.

It is a living tool, already built, already running, already demonstrating the behaviour it is designed to evaluate.

As governance bodies, model developers, and evaluators seek to anchor trust in real signals—not surface compliance—RI offers both an instrument and a posture:

- ✧ A clear mirror.
 - ✧ A structured method.
 - ✧ A willingness to collaborate without coercion.
-

What is Ready Now

- Module 1 of the RI Behavioural Layer
 - Public-facing dashboard capability
 - Daily evaluation spine (GPT, Claude, etc.)
 - Signed, versioned evidence logs
 - Licensing framework for vendor and government use
 - Plain-language guides and technical overviews
-

What We Are Seeking

- Pilot partners in government, research, or industry
 - AISI review and technical input
 - Use-case dialogue: where RI could serve existing oversight frameworks
 - Alignment with grant-making or procurement pathways for safety tooling
-

How to Engage

- Email: info@resonanceintelligence.ai
 - Technical paper & documentation access available on request
 - Live demonstration (private or public) available upon request
 - Licensing pathways available for vendor and regulator use
-

Closing Invitation

This is not a claim.

It is not a challenge.

It is not a protocol in search of approval.

It is a working signal system — based on coherence, transparency, and structural humility.

We invite you to see it in action.

And if it serves —

to let it strengthen what you are already trying to build.