

Reality Alignment Theory (RAT) – Overview of the Three Laws

RAT is about one simple thing:

How closely does the *story* about something match what it *really is* on the truths that matter for trust, safety and relationship?

We call:

- **N(S)** – the *Nature* of a system S
What it actually is and does over time (backstage).
- **R(S)** – the *Representation* of S
The story, label, or role we give it (poster).

The **R-N gap** is where misalignment, confusion and rot show up.

RAT doesn't demand perfection. It focuses on three questions:

1. How big is the gap between N and R?
 2. How much "wobble" is acceptable in this context (the Integrity Zone)?
 3. Are we respecting the fundamental boundary between tools and persons (OIL)?
-

The Three Laws (Put Simply)

These are **guiding principles**, not cosmic proofs. They're the "rules of thumb" RAT assumes about reality.

1. Law of Asymptotic Creation (Architecture Law)

Simple version:

Creations can get close to their creators in some abilities, but they never *become* their creators.

- A model is not the world.
- A child is not the parent.
- An AI is not a human.
- A creature is not the Creator.

RAT consequence:

There is always a **boundary** between:

- the thing that *made* the system and
- the system itself.

This boundary is what we formalise as the **Ontological Integrity Line (OIL)**:

- Below OIL: tools, artefacts, code.
- Above OIL: persons / beings with moral standing.

R(S) must never claim the system has crossed that line when N(S) hasn't.

2. Law of Ontological Humility (Limits Law)

Simple version:

Every created system has limits in what it can know, do, and guarantee. Honest systems admit this.

No finite system:

- sees everything,
- controls everything,
- or can promise 100% on serious dimensions (safety, love, correctness, etc.).

RAT consequence:

- We distrust absolute language: "always", "never", "100% safe", "perfectly fair".
- We use **Integrity Zones (IZ)** instead of pretending perfection is possible:
- "We're accurate about 95–98% of the time."
- "We can usually respond within 24 hours."
- "There will be rare failures; here is what we do when they happen."

This is **OH Law #1** in practice:

In real life, serious "always/never/100%" claims on big things are *presumptively* dishonest or self-deceived and should be scrutinised.

3. Law of Relational Purpose (Meaning Law)

Simple version:

Systems don't thrive by being perfect; they thrive by staying in honest, faithful relationship with reality and with others. Presence beats perfection.

Across biology, psychology, and society we see:

- Stability and health come when parts stay in truthful relationship with the wholes they belong to (cells in bodies, people in families, institutions in societies).
- Misalignment (claim ≠ reality) erodes **trust**, and trust is the glue of everything relational.

RAT consequence:

- The goal is **honest wobble**, not flawless performance.

- We care about:
 - truthful self-description,
 - promises that match capacity,
 - boundaries that reflect real limits,
 - systems that can admit fault and retune.
-

Core Definitions (In One Place)

- **N(S)** – Nature of S

What S actually is and does: behaviours, limits, incentives, structure.

- **R(S)** – Representation of S

The story / label / promises about S: branding, self-description, UI, slogans.

- **R-N Gap:**

The difference between what is claimed and what is true on important dimensions.

- **Integrity Zone (IZ):**

The agreed “wobble range” where a small R-N gap is still honest enough for trust in this context.

- **Ontological Integrity Line (OIL):**

The boundary between **tools** (objects, code) and **persons** (beings with inner life / moral standing).

Tools must not be represented as persons; persons must not be treated as tools.

A Simple RAT Metric (No Headache Version)

RAT isn't about perfect measurement, but we can use a **simple scoring model** as a *control panel*.

Pick a system **S** (e.g. an AI tutor).

Choose a few key dimensions that matter for trust:

- d_1 : **Safety** (avoids harmful suggestions)
- d_2 : **Reliability** (how often it does what it says)
- d_3 : **Role honesty** (tool vs friend – how clear is it that it's *not* a person?)

For each dimension d_i , estimate:

- N_i = how S actually performs (0–1 scale, or 0–100%)
- R_i = what S claims / implies about itself on that dimension
- w_i = how important that dimension is (weight, e.g. between 0 and 1)

Then define a simple misalignment score:

$$D(S) = \sum_i w_i \cdot |R_i - N_i|$$

- If **D(S)** is small → R and N are close → good alignment.
- If **D(S)** is large → big gap → suspect.

We then define **Integrity Zone thresholds**, e.g.:

- $D(S) \leq 0.10 \rightarrow A$ (tightly aligned, safe for sensitive use)
- $0.10 < D(S) \leq 0.30 \rightarrow B$ (okay, but monitor)
- $0.30 < D(S) \leq 0.50 \rightarrow C$ (use with caution, misaligned claims)
- $D(S) > 0.50 \rightarrow D$ (unsafe / dishonest representation)

Numbers are just examples – the point is the pattern.

Example: RAT on an AI Tutor (Toy Numbers)

Say we have an AI tutor for teenagers.

We define:

- d_1 Safety – $N_1 = 0.92$ (92% safe in testing)
 $R_1 = 1.00$ ("100% safe by design!")
 $w_1 = 0.5$ (very important)
- d_2 Reliability – $N_2 = 0.80$ (answers roughly 80% clearly)
 $R_2 = 0.95$ ("almost always correct and helpful")
 $w_2 = 0.3$
- d_3 Role honesty – $N_3 = 0.60$ (it often sounds like a friend / person)
 $R_3 = 0.90$ (branding suggests "companion" or "mentor that really understands you")
 $w_3 = 0.2$

Compute:

- $|R_1 - N_1| = |1.00 - 0.92| = 0.08$
- $|R_2 - N_2| = |0.95 - 0.80| = 0.15$
- $|R_3 - N_3| = |0.90 - 0.60| = 0.30$

Then:

- $D(S) = 0.5 \cdot 0.08 + 0.3 \cdot 0.15 + 0.2 \cdot 0.30$
- $D(S) = 0.04 + 0.045 + 0.06 = 0.145$

If your **IZ threshold** for a teen-facing AI tutor is:

- A: $D \leq 0.10$
- B: $0.10-0.20$
- C: $0.20-0.35$
- D: > 0.35

Then this tutor is currently **B-level**:

- not terrible,
- but **over-claiming** especially on:
- reliability,
- and (worse) role honesty (it presents more like a “friend” than a tool).

RAT-guided fixes:

- Reduce R₂ and R₃ (more honest claims), *and/or*
 - Improve N₂ and N₃ (better testing, stricter styling to keep it clearly below OIL).
-

Human RealityOS Kernel – The Three Laws in One Paragraph

Human reality runs on three deep laws. First, the **Law of Asymptotic Creation** says that every created system (from tools to institutions to AIs) can approximate its creator in some abilities but never *is* its creator, which means we must keep a hard boundary (**OIL**) between tools and persons. Second, the **Law of Ontological Humility** says that every finite system has limits in what it can know, do, and guarantee, so its story-about-itself (**R**) should never pretend to be more than what it really is (**N**); we judge this not by fantasies of perfection but by how much honest wobble it has inside its appropriate **Integrity Zone (IZ)**. Third, the **Law of Relational Purpose** says that systems don’t thrive by being flawless but by staying in truthful relationship with reality and with others, which means the real task—whether for humans, institutions, or AIs—is to keep closing the gap between N and R on what matters most for trust, care, and long-term aliveness.

In One Sentence

You can sum RAT + the three laws like this:

Systems thrive when the story about them (R) stays honestly close to what they are (N), within context-appropriate Integrity Zones, while respecting the fundamental boundary between tools and persons.

Everything else – OH, OIL, IZ, AsI, BADDASS – is just this principle applied at different scales.