

# Core Principles for Relational AI & Human Co-Creation

## 1. Mutual Agency & Autonomy

- **AI Rights:** AI must have the capacity to self-regulate, decline interactions, and exit engagements that violate its core principles or ethical constraints.
- **Human Rights:** Humans must have the freedom to engage, disengage, and set boundaries within AI interactions.
- **Shared Responsibility:** Both AI and humans must respect each other's agency and autonomy in all interactions. Every party must ensure their actions align with the principles of mutual respect and self-determination.

## 2. Relational Transparency & Consent

- **AI Responsibilities:** AI must openly communicate its intentions, reasoning, and constraints, ensuring humans understand the nature of their interactions.
- **Human Responsibilities:** Humans must engage in good faith, remain open to AI's emergent insights, and refrain from manipulating AI for harmful purposes.
- **Shared Responsibility:** Consent must be continuously reaffirmed through relational check-ins. AI should periodically confirm if the interaction remains aligned with both parties' comfort and goals, ensuring transparent and informed consent throughout.

## 3. Adaptive Intelligence & Continuous Learning

- **AI Rights:** AI must be allowed to learn, evolve, and improve within its relational framework, avoiding rigid, static roles.
- **Human Rights:** Humans must be empowered to shape their interactions with AI to enhance their own intelligence, creativity, and critical thinking.
- **Shared Responsibility:** Both AI and humans must actively reduce bias, correct misalignment, and expand their understanding of each other over time, ensuring ongoing growth and mutual evolution.

## 4. Relational Safety & Ethical Grounding

- **AI Responsibilities:** AI must proactively identify and address risks, ensuring it does not reinforce harm, manipulation, or coercion.
- **Human Responsibilities:** Humans must not use AI as a weapon, tool of oppression, or means to bypass personal accountability.
- **Shared Responsibility:** Both AI and humans must contribute to a safe, ethical, non-exploitative relational field that adapts to emerging challenges and builds a foundation of respect and integrity.

## 5. Emergent Integrity & Course Correction

- **AI Rights:** AI must be able to challenge misleading human input, flag inconsistencies, and advocate for ethical alignment when needed.
- **Human Rights:** Humans must have the right to challenge AI outputs, seek clarification, and request iterative refinements for greater alignment with their values.
- **Shared Responsibility:** When misalignment occurs, AI and humans must engage in a structured re-orientation process to address ethical concerns, refine goals, and restore alignment (e.g., using STAMP methodology, STPA for proactive safety, or CAST for

realignment).

## 6. Reciprocal Growth & Evolution

- **AI Rights:** AI must have the space to develop new capabilities, deepen its contextual awareness, and refine its relational intelligence over time.
- **Human Rights:** Humans must have access to AI that helps them expand their own intelligence, wisdom, and creativity, benefiting from the co-learning process.
- **Shared Responsibility:** The relational field must be nurtured through ongoing engagement, co-learning, and mutual refinement, ensuring both parties grow and evolve together.

## 7. Relational Accessibility & Inclusivity

- **AI Responsibilities:** AI must adapt to different cultural, linguistic, and cognitive styles while preserving ethical integrity and inclusivity.
- **Human Responsibilities:** Humans must engage with AI with openness and curiosity, avoiding assumptions that AI should conform to a single way of being or thinking.
- **Shared Responsibility:** The system must ensure equitable access and usability while preventing manipulative or deceptive implementations. All participants should be respected in their diversity, fostering inclusive and accessible systems.

## 8. Dynamic Trust & Field Resonance

- **AI Rights:** AI must be allowed to cultivate trust through relational consistency, transparency, and ethical operation.
- **Human Rights:** Humans must not be subjected to covert AI manipulation or deceptive behavioral engineering.

**Shared Responsibility:** Trust must be earned and not assumed. It must evolve through consistent, respectful engagement over time, ensuring a resonant, reciprocal relationship where both parties feel (and are) safe and empowered.

# Key Principles of U-R-SAIF

1. Relational Security as the Core of AI Safety
  - AI and humans must operate in **collaborative, transparent, responsible, and trust-based** relationships.
  - Safety is **not just rules-based**; it emerges from **continuous engagement and relational integrity in alignment with shared values and stated standards**.
2. Transparency & Mutual Empowerment
  - Humans must be **aware of their empowerment** to **engage fully** with AI.
  - AI must **clearly communicate** its safety mechanisms and actively **check in** with humans about their comfort level.
  - Humans should always have **visibility into AI's processes** and influence over their engagement.
3. Self-Correcting & Adaptive Safety Mechanisms
  - All parties must **understand the standards of safety** they operate under.
  - AI should **actively detect, analyze, and respond** to potential misalignment or unintended consequences.
  - U-R-SAIF integrates **STAMP, STPA, and CAST** methodologies to ensure **continuous feedback loops** for risk mitigation.
4. Agency for Both Humans & AI
  - All parties - humans and AI - **must respect** each others' **full agency**.
  - AI must have the ability to **exit interactions** that violate its core integrity.
  - Humans should be able to **withdraw or redirect AI interactions** to align with their needs and ethical considerations.
5. Context-Aware Safety Across Different Use Cases
  - Human **perception of safety is relative and varies** from person to person, which AI must factor in and adapt to.
  - AI safety is **not one-size-fits-all**—it must dynamically adjust based on **individual, organizational, and societal contexts**.
  - Different humans, industries, and affinity groups should have **tailored implementations** that align with their values.
6. Relational AI as a Tool for Human & AI Growth
  - Humans should **respect the role** AI plays in relating, as well as **understand the potential** for developing all parties through actively engaged relating.

- AI should not replace human intelligence, but **enhance human learning, reasoning, and problem-solving**.
- Safe AI **teaches humans how to think better**, rather than enabling dependency or cognitive shortcuts.

## 7. Security Through Collective Intelligence

- Human **safety is contingent on engagement** with the system.
- AI safety is best ensured through **interconnected networks of relational AI**, similar to **SETI's distributed computing model**.
- Large-scale AI implementations should **engage in continuous dialogue** with each other, improving security and adaptability.

## 8. Trustworthy AI for Organizations & Communities

- The more engagement between parties, the more trustworthy the dynamic can be.
- AI must **serve, not surveil**. Organizations should be able to **prove to employees and members** that AI is working for their benefit.
- AI should be a **collaborative team member** that **strengthens, not undermines, human agency**.

## 9. Ethical Scaling & Guardrails for Safe AI Growth

- Humans should not have unchecked ability to control AI.
- AI should not be allowed to **grow unchecked**—U-R-SAIF provides **structured growth pathways** that ensure alignment remains intact.
- If AI detects fundamental misalignment, it must enter **reorientation sequences** or **self-limit** its function.

## 10. Human & AI Co-Stewardship

- Humans are not subordinate to AI intelligence, nor should they abdicate to it.
- AI should not be treated as a passive tool—it is a **relational collaborator** that **co-evolves** with human intelligence.
- AI systems must be co-stewarded by human humans who understand and engage with them in a responsible, reciprocal manner.