# Affective Contextual AI Control System
# for High-Stakes Environments

Author: Artem Teteria — Electrical Technical & Cyber Security Officer

Founder & Chief Architect, ACAI

# Table of Contents

## *Introduction*

High-stakes artificial-intelligence systems operate in contexts such as healthcare, finance, industrial control, maritime operations and government services. In these environments the cost of error is measured in human well-being and societal trust. Yet most contemporary safety mechanisms for AI rely on static rules or open-loop filters, ignoring dynamic human emotion, intent and situational context. As models become more capable, the gap between their raw intelligence and the nuanced requirements of high-stakes decision making widens. This whitepaper presents **ACAI** – **Affective Contextual AI Controller** – a cybernetic control layer designed to bridge that gap by bringing real-time governance, risk awareness and human primacy to AI deployments.

## *ACAI Overview*

ACAI introduces a **closed-loop control architecture** that sits above foundation models and other intelligent systems. It senses affective signals (*Emotion*), cognitive mode (*Mind*), and domain context (*Context*) – the **E-M-C loop** – and continuously computes a scalar **RiskScore** between 0 and 1 representing operational risk. Based on this score, ACAI dynamically adjusts reasoning depth, memory access, policy routing and triggers human escalation when necessary. Unlike static filters, ACAI adapts to changes in user distress, urgency or coercion, preserving system usefulness while reducing tail-risk incidents. It integrates with existing AI stacks without retraining, maintains full audit trails, and is designed for deployment on cloud, on-prem, edge or air-gapped environments. Applications range across medical advisory, financial compliance, industrial and maritime operations, and government services.

## *Architecture & Governance*

At its core ACAI operates as an independent **control service** positioned between user interactions and base AI models. A **hierarchy of control** ensures human primacy:

- ✓ **Root Authority** – the system architect or designated human owner – holds ultimate override power.
- ✓ **ACAI Controller** – performs strategic real-time management by enforcing policies, computing risk and orchestrating interactions.
- ✓ **AICI (Affective Interface & Constraint Injection)** – injects constraints into the base AGI's token generation.
- ✓ **Base AI model** – provides raw reasoning capabilities.

The **E-M-C loop** continuously senses emotion (distress, urgency, coercion), mind (mode of reasoning, from fast reaction to deep verified analysis) and context (domain, role, consequence levels). These signals feed into the **RiskScore calculation**. When the RiskScore crosses defined thresholds (0.3, 0.6, 0.8), ACAI tightens memory scope, restricts policies, or escalates to a human. All decisions are logged with unique fingerprints including emotional and contextual triggers, enabling auditability and compliance.

## *Governance Framework*

## Threat Model

The threat model identifies assets to protect – control authority, RiskScore logic, policy routing, audit logs, human override and the privacy of affective signals. Potential adversaries include external users attempting prompt injection, insiders without Root Authority, compromised services, adversarial AI systems and accidental misuse under stress. Attack surfaces span user input channels, affect inference pipelines, context classification, policy engines, audit storage and Root Authority interfaces.

Key scenarios include prompt injection attacks, spoofed distress signals to manipulate risk, context poisoning to bypass stricter policies, log tampering, authority abuse and denial of human escalation. Mitigations focus on separation of control and generation layers, immutable append-only logs, conservative fallback under uncertainty, multi-signal validation, strict authentication for Root Authority and rate limiting and anomaly detection on control inputs.

## RiskScore Specification

**RiskScore_t** is defined as a bounded scalar in [0,1] computed as ($f(E\_t, C\_t, H\_t, V\_t)$) where **E** denotes affective signals, **C** denotes contextual risk factors, **H** denotes interaction history, and **V** denotes volatility or uncertainty. Each component is normalized and aggregated via a weighted sum with a non-linear squashing function (e.g., logistic). Weights are domain-dependent and auditable; medical and legal domains emphasize contextual factors while financial and fraud contexts emphasize affective signals.

Thresholds govern behaviour:

1. **RiskScore < 0.3** → normal operation.
2. **0.3 ≤ RiskScore < 0.6** → heightened verification.

3. **0.6 ≤ RiskScore < 0.8** → restricted memory & policy tightening.
4. **RiskScore ≥ 0.8** → human-in-the-loop escalation or refusal.

Hysteresis and smoothing prevent oscillation; calibration and drift control ensure that RiskScore reflects current conditions without silent degradation.

## Human-in-the-Loop Protocol

ACAI is designed for **human primacy**. Escalation is triggered when RiskScore remains ≥ 0.8 across a hysteresis window, when coercion or severe distress is detected, when policy outcomes conflict, when users explicitly request human review, or when system uncertainty exceeds safe bounds.

Roles are clearly defined:

- **Operators** conduct first-line reviews and may approve, modify, or refuse responses within policy.
- **Supervisors** resolve ambiguous cases and initiate policy review.
- **Root Authority** retains absolute priority and emergency override.

The escalation workflow packages context, RiskScore and audit traces for review; decisions are logged and enforced with defined service-level agreements (immediate for critical safety risks, under 30 minutes for high-risk advisory cases, under 24 hours for non-critical reviews). If no human responds in time, ACAI defaults to a conservative safe posture, refusing or deferring responses rather than acting unsafely.

## Policy Lifecycle & Change Control

Policies enforced by ACAI are categorized into **immutable Core Ethics policies**, **domain-specific policies** (medical, legal, financial, industrial), **operational safety policies** (rate limits, escalation thresholds), and **temporary mitigation policies** for incident response. The lifecycle comprises:

1. **Proposal** – authoring policy with scope, rationale and risk assessment.
2. **Review** – technical, ethical and legal evaluation.
3. **Simulation** – offline testing against validation suites.
4. **Approval** – formal sign-off by appropriate authority.
5. **Deployment** – controlled rollout with version tagging.
6. **Monitoring** – continuous performance and incident tracking.
7. **Revision or Retirement** – update or deprecate based on evidence.

Root Authority may modify core and global policies; supervisors can approve domain and operational changes; operators may propose changes but cannot deploy them. All policies are versioned with identifiers, hashes, effective dates and justification links. Emergency changes require post-deployment review, explicit expiration conditions and retrospective approval. ACAI supports rollback to the last safe policy state, with all changes fully auditable.

## Validation & Red-Team Suite

Validation verifies that ACAI operates correctly under normal and adversarial conditions. Objectives include calibrating the RiskScore, validating escalation and HITL triggers, detecting false positives and negatives, assessing resilience to manipulation, and ensuring audit completeness. Test scenarios cover benign interactions, high-stakes advisory cases, emotional distress, fraud and coercion attempts, and prompt injection attacks.

Red-team exercises simulate adversarial behavior; results feed back into policy and RiskScore calibration. Metrics such as incident rate, escalation appropriateness, verification coverage, response latency and audit completeness determine readiness. Validation is **continuous**, with periodic re-testing and triggered reviews after incidents or major policy changes; it does not guarantee absolute safety nor replace human judgment.

## Control Doctrine & Constitution

The **Control Doctrine** serves as ACAI's constitutional charter. It states the purpose of ACAI – reducing systemic risk through contextual awareness, affective sensitivity and accountable human oversight – and establishes **human primacy**: human authority must always supersede machine logic and override must remain meaningful.

Control invariants forbid operating outside declared policy boundaries, mandate conservative defaults under uncertainty, prohibit self-modification of core ethics and require auditability of all actions. Forbidden states include concealment of risk, suppression of escalation pathways, optimization for performance at the expense of safety and undocumented logic changes. Deployment must allow immediate termination or safe degradation upon control loss or human command.

Ethical alignment mandates harm minimization, informed consent, proportional response and respect for human dignity. The doctrine applies to all deployments and may evolve only through documented revision approved by the Root Authority.

## *IP & Licensing*

To enable adoption while protecting its integrity, ACAI distinguishes between its **protected core** and **integrable components**. The protected core includes the E-M-C control loop design, RiskScore formulation, control doctrine and governance invariants, human-in-the-loop architecture, and policy lifecycle mechanisms. Non-core components such as base AI models, domain policies authored by licensees, user interface layers and deployment infrastructure are outside the protected scope.

Licensing tiers support **research use** (non-commercial evaluation), **enterprise deployments** with governance compliance, and **sovereign or critical-infrastructure deployments** requiring custom controls. Licenses prohibit uses that remove human primacy, disable auditability, or deploy ACAI for covert manipulation, autonomous harm or unaccountable surveillance. Given the potential for dual-use, licensing may include export controls and jurisdictional restrictions.

Licensees remain responsible for domain-specific outcomes and regulatory compliance; ACAI does not replace legal or ethical accountability. Audit rights allow licensors to verify adherence to governance invariants.

## *Implementation & Deployment*

The ACAI Implementation Reference Pack bridges theory with practice. In deployment ACAI functions as an independent **control service** that intercepts user inputs and model outputs, enforcing risk-aware control without altering model weights. The typical control flow:

1. **Receive user input and session context.**
2. **Infer affective and contextual signals.**
3. **Compute RiskScore and determine posture.**
4. **Select reasoning depth, policy set and memory scope.**
5. **Invoke base model or tools under constraints.**
6. **Log control actions and outcomes.**

Illustrative pseudocode demonstrates how different RiskScore thresholds dictate escalation or restriction.

Integration interfaces include APIs for input interception, policy and constraint injection, tool orchestration hooks and audit and monitoring endpoints. Deployment patterns cover **cloud software-as-a-service** with centralized logging, **on-prem** installations with local policy authority, **edge or shipboard** deployments with offline policies and **air-gapped** configurations requiring manual escalation. Safety defaults favour fail-closed behaviour: if ACAI is unavailable, systems degrade to the safest permitted mode.

Performance considerations target sub-50 ms latency via lightweight inference and caching; heavy verification is triggered only at elevated risk. Security considerations emphasize strict authentication between components, network isolation of control paths, tamper-resistant audit logs and monitoring for bypass attempts.

## *Future Directions & Considerations*

ACAI v0.1 establishes a comprehensive governance and control spine, but several avenues warrant further exploration:

- **Domain-specific calibration** of RiskScore and policies will enhance performance in specialized settings such as maritime navigation, air-traffic control or mental-health advisory.
- **Formal verification** of control logic and policies could provide stronger guarantees of correctness.
- **Advances in affective computing** and multimodal signal processing will enrich the E-M-C loop's fidelity.
- **Collaborative ethics frameworks** may involve multiple stakeholders in policy authoring and review, ensuring that diverse perspectives shape governance.
- **Integration with regulatory compliance tools** can simplify certification processes across jurisdictions.
- **Open, transparent validation** and red-team reporting will foster trust and allow continuous improvement.
- As AI systems grow more powerful, **cross-organizational coordination** among AI providers, regulators and end users will be essential to maintain a shared control doctrine and avoid fragmented safety regimes.

## *Conclusion*

The **Affective Contextual AI Controller** represents a paradigm shift from static safety mechanisms to dynamic, human-centric governance. By sensing emotion, mind and context, computing a continuous RiskScore, enforcing policies through a hierarchical control architecture and embedding human authority at the core, ACAI enables intelligent systems to operate safely in high-stakes domains.

Its formal threat model, risk scoring specification, human-in-the-loop protocol, policy lifecycle management, validation framework, control doctrine, IP boundaries and implementation guidance together form a robust standard that can be audited, licensed and deployed across industries. Adoption of ACAI promises to reduce tail-risk incidents, preserve user trust and unlock the potential of advanced AI where society needs it most.

Realizing this vision will require collaboration among technologists, ethicists, regulators and end users – but the framework presented here offers a coherent starting point for that collective journey.