System Design for Ethical Safeguards

Embedding accountability, fairness, and transparency directly into the architectural core (KA-MOD-013).

The Three Pillars of Ethical Integration



Proactive Assessment

Identify Risks Before Development

- **Bias Audit:**
 Analyze training data for demographic skew.
- **Impact Mapping:**
 Determine potential
 adverse effects on
 users/groups.
- **Privacy by Design (PbD):** Integrate data minimization from the start.



Architectural Encoding

Hard-Code Protections into the System

- **Separation of Duties:** Enforce different access levels for critical functions.
- **Explainability (XAI):** Design system logs to track and justify decisions.
- **Differential
 Privacy:** Add noise to data queries to protect individual records.



Continuous Operation

Monitor and Govern the Live System

- **Red Teaming:** Continuous attempts to exploit ethical vulnerabilities.
- **User Feedback Loops:**
 Dedicated channels for reporting perceived bias or harm.
- **Automated Auditing:**
 Scheduled checks for compliance drift and data integrity.

Key Design Requirements



Non-Repudiation

A strong, immutable audit trail proves who accessed what and when.



Fairness Metrics

Establish and monitor specific metrics to detect and mitigate algorithmic bias.



Reversibility

The ability to retract or undo a critical action if an ethical violation is detected.



Transparency

Documentation and clear interfaces explaining system logic to authorized users.