

Robert McCoy

Responsible Applications of AI

AIML-510-01

---

**Responsible AI Governance Framework**

---

“Making AI Understandable, Ethical, and Actionable.”

---

Mission: "To guide high-stakes organizations through responsible AI deployment, ensuring systems are ethically aligned and operationally sound."

**RM AI**  
**PROJECT CONSULTING**

# Ethical Principles Statement

**Fairness & Equity** – AI systems must avoid bias and ensure decisions are just across all users and contexts.

**Accountability & Oversight** – Every deployment has a human owner responsible for outcomes, supported by governance review boards.

**Transparency & Explainability** – Models are documented, explainable, and understandable by stakeholders.

**Privacy & Protection** – User data is safeguarded under GDPR, CCPA, and evolving compliance regimes.

**Sustainability & Long-Term Stewardship** – AI adoption considers societal, workforce, and environmental impacts.

# RAI Governance & Oversight Model

- **Executive Oversight** – Senior leadership ensures AI strategy aligns with mission and ethics.
- **RAI Council / Board** – Cross-functional group (Legal, IT, Product, Ethics, HR) approves major deployments.
- **Operational Teams** – Developers, data scientists, and compliance staff apply governance policies in daily workflows.
- **External Advisors** – Periodic independent audits for transparency and accountability..

# Roles & Responsibilities

- **Chief AI Ethics Officer** – Chairs RAI Council; ensures fairness, bias audits, and ethical alignment
- **Legal & Compliance Lead** – Interprets GDPR, CCPA, and EU AI Act; integrates regulatory requirements into AI contracts
- **Data Science & Engineering Teams** – Apply bias testing, document model cards, maintain explainability tools
- **IT & Security Lead** – Safeguards data integrity, access control, and SOC2/FedRAMP compliance
- **HR & Training Lead** – Develops employee RAI literacy programs, conducts workforce awareness sessions.
- **End Users / Business Owners** – Final accountability for AI outcomes; empowered through monitoring and training.

# Standing Up Responsible AI – Phased Roadmap



# Measuring Responsible AI Success

## Metric Area

### Fairness & Bias

## Example Measures

- % of bias tests passed (NIST SP 1270 aligned)
- Number of fairness audits completed per quarter

### Transparency & Explainability

- % of models with published model cards / datasheets
- Average stakeholder explainability rating (survey)

### Privacy & Security

- Compliance with GDPR/CCPA/FedRAMP
- Number of privacy incidents or breaches reported

### Accountability & Oversight

- % of projects reviewed by RAI Council
- Time to resolution of flagged ethical issues

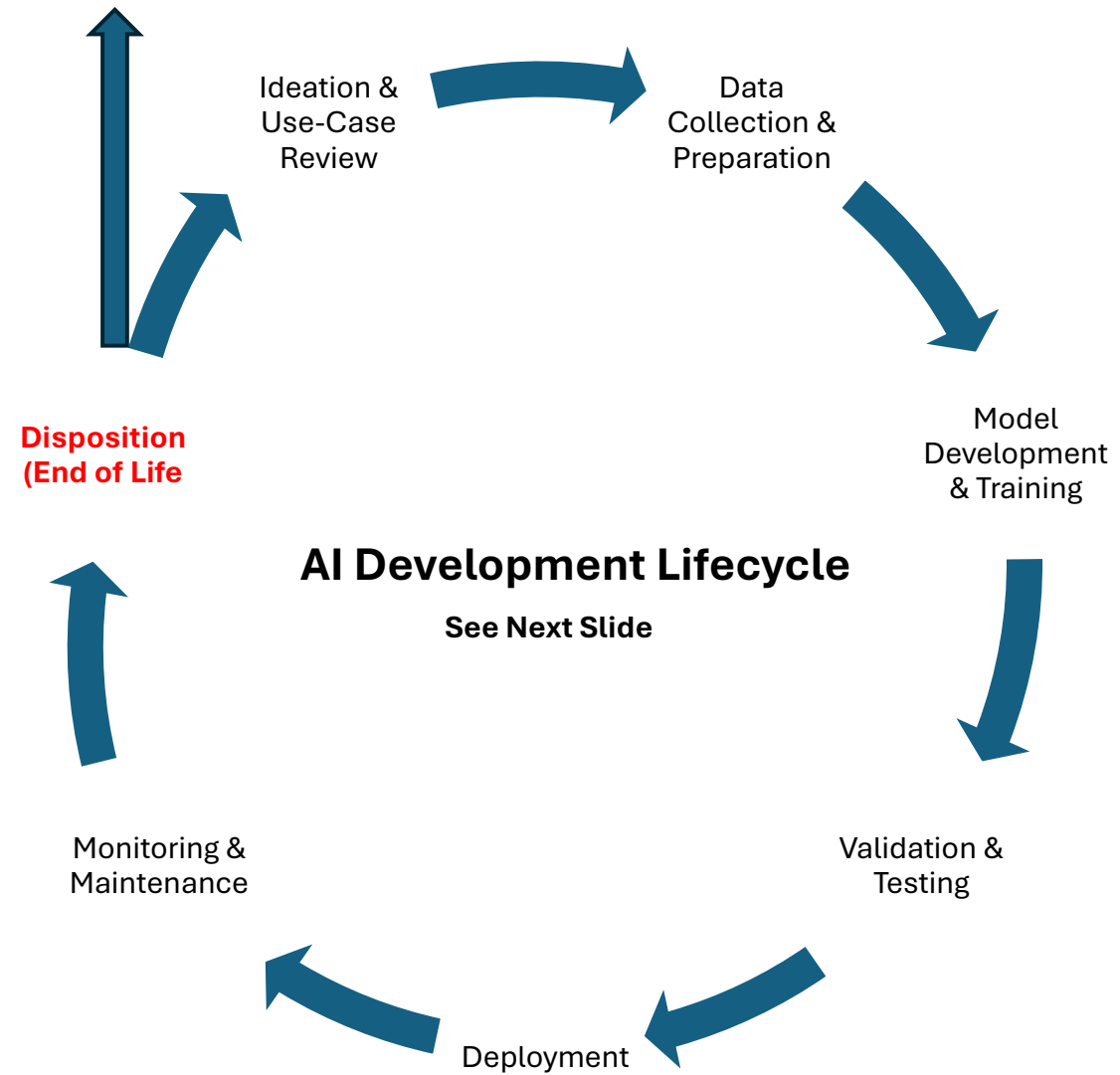
### User Trust & Training

- Employee RAI training completion rate
- User trust scores from post-deployment surveys

### Sustainability & Long-Term Impact

- Reduction in skill atrophy risk (tracked via user monitoring)
- Evidence of positive workforce/customer impact

# RAI Across the AI Development Lifecycle





# RAI Across the AI Development Lifecycle

## Lifecycle Stage

### Ideation & Use-Case Review

### Data Collection & Preparation

### Model Development & Training

### Validation & Testing

### Deployment

### Monitoring & Maintenance

### Disposition (End-of-Life)

## RAI Actions & Controls

- Assess ethical implications & intended benefits
- Screen for high-risk applications (e.g., medical, defense)
- Apply bias taxonomy (systemic, statistical, cognitive)
- Ensure privacy compliance (GDPR, CCPA)
- Document data lineage & governance
- Perform fairness & robustness testing
- Use diverse datasets & augmentation strategies
- Maintain model cards & transparency records
- Run red-team/blue-team simulations
- Conduct explainability reviews
- Verify compliance with internal RAI scorecards
- Secure access & SOC2/FedRAMP readiness
- Require human-in-the-loop for critical decisions
- Publish user-facing explainability statements
- Track bias drift, accuracy decay, and fairness gaps
- Continuous user feedback integration
- Scheduled RAI audits & incident reporting
- Decommission responsibly (data purged, access revoked)
- Archive documentation for audit trail
- Ensure no “orphaned” models remain active



# Transparency & Explainability

## Transparency Dimension

### Model Documentation

### Explainability Tools

### Stakeholder Communication

### Auditability

### Balancing Transparency & IP

## RAI Actions & Best Practices

- Publish model cards & datasheets • Include purpose, limitations, and known risks
- Apply SHAP, LIME, or equivalent interpretability methods • Provide plain-language outputs for stakeholders
- Create user-facing “explain it simply” summaries • Run fairness workshops with impacted groups
- Maintain governance-ready documentation • Log all significant model decisions for external/internal review
- Share enough detail for trust without exposing trade secrets • Use tiered disclosure (public vs. regulator vs. internal)

# Identifying and Mitigating Bias in AI Systems

## Bias Source

**Systemic Bias** (historic inequities baked into data)

**Statistical Bias** (imbalanced or skewed datasets)

**Human / Cognitive Bias** (annotation errors, subjective labels)

**Evolving Bias** (shifts in norms & environment over time)

**Lack of Ground Truth** (flawed human data, incomplete truth sets)

## RAI Mitigation Actions

- Use participatory design reviews
- Diversify data collection across demographics
- Apply fairness constraints during model training
- Apply reweighting & resampling techniques
- Use synthetic/augmented data
- Validate with independent test sets
- Rotate and train annotators
- Blind labeling where possible
- Consensus validation on disputed cases
- Ongoing monitoring & bias bounties
- Regular model retraining
- Scenario testing with stakeholder groups
- Document uncertainty transparently
- Avoid overclaiming model accuracy
- Incorporate external audits of fairness

# Privacy, Protection & Compliance

## Focus Area

### Data Privacy

### Data Protection

### Regulatory Compliance (U.S.)

### Regulatory Compliance (EU)

### Intellectual Property & Ownership

## RAI Practices & Controls

- Apply data minimization & purpose limitation
- Anonymize or pseudonymize sensitive data
- Obtain clear, informed consent from users

- Encryption at rest and in transit
- Strict access controls & role-based permissions
- Regular penetration testing and audits

- Follow CCPA, HIPAA, and sector-specific rules
- Maintain SOC2 compliance
- Document consent and data usage policies

- Align with GDPR and EU AI Act
- Ensure right-to-explanation and portability
- Conduct Data Protection Impact Assessments (DPIAs)

- Clarify IP rights when training on mixed datasets
- Contract clauses for explainability, fairness, and liability
- Define ownership of AI-generated outputs

# User Monitoring, Training & Education

## Focus Area

### User Monitoring

## RAI Actions & Practices

- Track AI usage patterns for anomalies
- Detect over-reliance or automation bias
- Establish escalation protocols for misuse

### Employee Training

- Mandatory onboarding modules on bias, transparency, and compliance
- Annual refresher courses with updated regulations
- Scenario-based simulations (ethical dilemmas, misuse cases)

### Customer Education

- Provide clear explanations of AI features & limitations
- Offer “AI 101” guides and FAQs
- Transparent consent and opt-out options

### Feedback & Engagement

- User surveys on trust and satisfaction
- Fairness workshops with impacted groups
- Continuous improvement from user feedback loops

### Skill Preservation

- Encourage critical thinking alongside AI
- Promote human-in-the-loop workflows
- Monitor for skill atrophy risk (e.g., medical, aviation)

# Considering Long-Term Effects on Employees and Customers

## Impact Area

### Workforce & Skills

### Employee Well-Being

### Customer Trust

### Societal & Ethical Impact

### Sustainability & Environment

## RAI Considerations & Actions

- Guard against skill atrophy from over-reliance
- Reskill employees for AI-augmented roles
- Support career transition programs
- Monitor stress linked to constant AI oversight
- Establish clear boundaries on surveillance
- Foster a culture of ethical empowerment
- Build transparency into every product touchpoint
- Provide clear channels for dispute resolution
- Prioritize explainability in consumer-facing AI
- Mitigate systemic bias and inequities over time
- Ensure AI reflects aspirational ethics, not just current norms
- Support community dialogue and public education
- Evaluate AI's energy consumption
- Favor efficient architectures & green compute practices
- Align adoption with corporate sustainability goals

# Best Practices & Open Questions

## Best Practice

**Cross-Functional Governance** – Involve Legal, IT, Ethics, HR, and Product in AI oversight

**Bias Audits & Fairness Reviews**  
Conduct regular audits with NIST SP 1270 alignment

**Transparent Documentation** – Maintain model cards, datasheets, and bias mitigation logs

**User-Centric Training** – Equip employees and customers to use AI responsibly

**Independent Audits** – Engage external advisors for trust and accountability

## Open Question / Challenge

Who ultimately “owns” Responsible AI in an organization?

Can “fairness” ever be universally defined across cultures?

How do we balance transparency with intellectual property protection?

How do we prevent long-term skill erosion from over-reliance?

What accountability exists if a system meets technical specs but causes social harm?



# References

- European Commission. (2021). *Proposal for a Regulation of the European Parliament and of the Council laying down harmonised rules on artificial intelligence (Artificial Intelligence Act)*. COM/2021/206 final. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A52021PC0206>
- European Union. (2016). *General Data Protection Regulation (GDPR)*. Regulation (EU) 2016/679. <https://eur-lex.europa.eu/eli/reg/2016/679/oj>
- Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2017). Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in Neural Information Processing Systems*, 30, 6402–6413. <https://arxiv.org/abs/1612.01474>
- NIST. (2023). *AI Risk Management Framework (NIST AI RMF 1.0)*. National Institute of Standards and Technology. <https://www.nist.gov/itl/ai-risk-management-framework>
- NIST. (2022). *Towards a Standard for Identifying and Managing Bias in Artificial Intelligence (NIST SP 1270)*. National Institute of Standards and Technology. <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.1270.pdf>
- Parasuraman, R., & Riley, V. (1997)., 39(2), 230–253. <https://doi.org/10.1518/001872097778543886>
- Shneiderman, B. (2020). Human-centered artificial intelligence: Reliable, safe & trustworthy. *International Journal of Human-Computer Interaction*, 36(6), 495–504. <https://arxiv.org/pdf/2002.04087>
- U.S. Congress. (2018). *California Consumer Privacy Act (CCPA)*. California Civil Code § 1798.100. [https://ccpa.ca.gov/regulations/pdf/ccpa\\_statute.pdf](https://ccpa.ca.gov/regulations/pdf/ccpa_statute.pdf)
- World Economic Forum. (2021). *Global AI Action Alliance: A framework for responsible AI adoption*. Geneva: World Economic Forum. <https://www.weforum.org/press/2021/01/world-economic-forum-launches-new-global-initiative-to-advance-the-promise-of-responsible-artificial-intelligence/>
- Zuboff, S. (2019). *The age of surveillance capitalism: The fight for a human future at the new frontier of power*. New York: PublicAffairs.