# Responsible AI Governance for Internal Audit

*What Internal Auditors Need to Know About AI Governance Standards, Controls, and Oversight Models*

## 1. Why Responsible AI Governance Matters

AI introduces new risks (accuracy, fairness, drift, misuse, security) that must be managed with **clear governance protocols, defined escalation paths, formal oversight, and transparent documentation**.

The goal of this handout is to provide:

- A clear summary of global AI frameworks
- Operational governance expectations
- Escalation paths and intervention requirements
- Control mapping considerations
- Fallback and fail-safe mechanisms
- Assurance and reporting expectations
- Alignment to ISO, NIST, OECD, EU, and IEEE standards

Designed specifically for **internal audit leaders and practitioners**.

## 2. The Five Pillars of Responsible AI Governance

1. **Structured oversight:** Policies, roles, responsibilities, and committees provide decision authority and accountability.
2. **Lifecycle governance:** Controls exist from model conception → design → development → deployment → monitoring → retirement.
3. **Human oversight and intervention:** Clear thresholds for when humans must step in and override the system.
4. **Transparency and documentation:** AI systems, inputs, outputs, assumptions, decision paths, and change history must be visible and auditable.
5. **Ethical and legal compliance:** Controls prevent harmful, biased, or unintended outcomes and maintain compliance with global regulations.

# 3. Governance Protocols (What Must Exist)

| | | | |
|---|---|---|---|
| A | **Escalation and Human Intervention Protocols** | *Every AI system should define*:<br>• When a human must intervene<br>• How exceptions are routed<br>• Who is authorized to override AI decisions<br>• What risk thresholds require escalation<br>• What scenarios trigger shutdown or fallback | *Internal audit should confirm*:<br>• Escalation paths are documented and tested<br>• Staff are trained on intervention criteria<br>• Oversight responsibilities are clearly assigned |
| B | **Control Mapping: Human vs. AI Responsibilities** | *AI governance requires clarity about*:<br>• What the AI system is allowed to do<br>• What decisions require a human<br>• How tasks shift between automation and manual work<br>• How control activities change in mixed human-AI environments | *IA should verify*:<br>• Updated RACI charts<br>• Clear human-in-the-loop or human-on-the-loop expectations<br>• Control redesign where AI replaces or augments tasks |
| C | **Fallback & Fail-Safe Mechanisms** | *When an AI model fails, behaves unpredictably, or produces low-confidence results, organizations must have*:<br>• Fallback rules<br>• Safe-mode behavior<br>• Manual alternatives<br>• Auto-disable conditions | *IA evaluates*:<br>• Whether fallback rules exist<br>• Whether they align with enterprise risk tolerance<br>• Whether they have been tested under real conditions |
| D | **Traceability, Logging & Assurance Reporting** | *All AI activity must be **visible and auditable**. Expectations include*:<br>• Logs of both human and AI decisions<br>• Model change histories<br>• Performance monitoring records<br>• Alerts, overrides, and exceptions | *Internal audit must confirm*:<br>• Logs are complete, retained, and reviewable<br>• Reporting captures both human and AI activity<br>• Evidence supports accountability |

## 4. The Major AI Governance Frameworks (Expanded)

Internal audit should use these standards to evaluate bias design controls; assess transparency artifacts; validate risk integration across functions; and benchmark auditor competency requirements.

| | Framework | Focus | Key Expectations | IA Relevance |
|---|---|---|---|---|
| **A** | **EU AI Act** | Risk classification + mandatory controls for high-risk AI | • Risk assessment<br>• Human oversight<br>• Technical documentation<br>• Monitoring and logging<br>• Data governance<br>• Vendor accountability | • Verify classification of AI systems<br>• Confirm compliance with oversight and logging requirements<br>• Review vendor documentation for high-risk systems |
| **B** | **ISO/IEC 42001** | Enterprise governance & lifecycle management | • AI governance structure<br>• Policies & procedures<br>• Training & competencies<br>• Lifecycle controls<br>• Monitoring and incident mgmt.<br>• Alignment with ISO 27001 & 31000 | • Assess maturity of the AI governance system<br>• Validate lifecycle documentation<br>• Confirm integration with security & ERM processes |
| **C** | **NIST AI Risk Management Framework** | AI risk identification, measurement, minimization | • Govern<br>• Map<br>• Measure<br>• Manage | • Validate risk identification activities<br>• Review metrics and performance evaluations<br>• Confirm governance structures align with NIST RMF expectations |
| **D** | **OECD AI Principles** | Human-centered and ethical AI | • Fairness<br>• Transparency<br>• Accountability<br>• Safety and robustness | • Use these as a values-based lens for evaluating ethical risk |

**E. Additional Important Standards:**

- **ISO/IEC 23894 — AI Risk Management:** Practical risk assessment guidance specific to AI design and deployment.
- **ISO/IEC 42005 — AI Transparency Requirements:** Defines the information organizations MUST be able to explain.
- **ISO/IEC 42006 — Competence of AI Auditors:** Specifies knowledge, skills, and qualifications for internal/external AI auditors.
- **ISO 31000 — Enterprise Risk Management:** Provides the risk integration backbone for AI risk reporting.
- **IEEE 7000 — Ethics in System Design:** Framework for addressing ethical considerations proactively.

## 5. Competency Model for AI Governance

Internal audit teams need competency in four areas:

**1. AI Foundations**

- Types of AI (predictive, generative, agentic)
- Human/AI decision boundaries
- Data lifecycle & model inputs/outputs

**2. AI Governance & Oversight**

- Explainability expectations
- Fairness, accountability
- Monitoring and drift controls
- Stakeholder alignment

**3. Regulatory & Ethical Awareness**

- Global regulations (EU, US, APAC)
- OECD & G7 principles
- Cross-border risk considerations
- Ethical risk identification

**4. AI Risk & Assurance Frameworks**

- NIST AI RMF
- ISO/IEC 42001 lifecycle
- ISO/IEC 23894 risk mgmt.
- ISO 27001 security
- Mapping controls to frameworks

## 6. What Internal Audit Should Evaluate in Any AI System

✔ **Governance:** Governance structure, oversight committee, policies, documentation.

✔ **Data Governance:** Lineage, quality checks, bias testing, access controls.

✔ **Model Development:** Requirements, assumptions, versioning, testing.

✔ **Model Validation:** Accuracy, fairness, robustness, limitations.

✔ **Human Oversight:** Intervention triggers, approvals, override authority.

✔ **Monitoring & Drift:** Performance indicators, alert thresholds, retraining.

✔ **Controls & Security:** Access, change management, vulnerability mgmt.

✔ **Auditability:** Traceability, logs, decision capture, documentation.

✔ **Vendor Management:** Transparency, SOC reports, disclosures, contractual controls.

## UNIFIED AI GOVERNANCE MATRIX

*(Governance Domain → Internal Audit Focus → Evidence → Framework Mapping)*

| Governance Domain | What IA Evaluates | Evidence Required | Framework Alignment |
|---|---|---|---|
| AI Inventory & Classification | Existence, completeness, risk tiering | Inventory, use case register | EU AI Act, ISO 42001, NIST Govern |
| Policies & Governance Structure | Policies, committees, roles, RACI | Policy documents, charters | ISO 42001, NIST Govern |
| Data Governance | Lineage, quality, bias controls | Lineage diagrams, DQ logs | ISO 42001, NIST Map/Measure, OECD Fairness |
| Model Development | Requirements, documentation, code controls | BRD, model card, version control logs | NIST Map, ISO 42001 lifecycle |
| Model Validation | Independence, accuracy, fairness | Validation reports, test results | EU AI Act, NIST Measure |
| Human Oversight | Intervention criteria, escalation | SOPs, override logs | EU AI Act, OECD Accountability |
| Monitoring & Drift | Alerts, KPIs, retraining | Drift dashboards, performance metrics | NIST Manage, ISO 42001 |
| Security & Access Control | Access reviews, change control | Access lists, change tickets | ISO 27001, ISO 42001 |
| Ethical & Compliance Risk | Bias, harm prevention, transparency | Bias test reports, transparency docs | OECD, IEEE 7000 |
| Vendor and 3rd-Party Governance | Transparency, assurances, contracts | SOC2, model cards, SLAs | EU AI Act, ISO 42001 |

## CROSS-FRAMEWORK HARMONIZATION MATRIX

*(Unifies EU AI Act, ISO 42001, NIST AI RMF, OECD, IEEE 7000)*

| Governance Pillar | EU AI Act | ISO 42001 | NIST AI RMF | OECD Principles | IEEE 7000 |
|---|---|---|---|---|---|
| Risk Classification | Strong | Moderate | Moderate | N/A | N/A |
| Transparency | High-risk transparency | High | High | High | High |
| Documentation | Required | Required | Expected | Minimal | Required |
| Human Oversight | Mandatory | Strong | Present | Strong | Strong |
| Data Governance | Required | Strong | Strong | Strong | Moderate |
| Monitoring & Drift | Required | Required | Required | Weak | Weak |
| Ethical Controls | Limited | Strong | Strong | Strong | Very Strong |
| Lifecycle Governance | Moderate | Very Strong | Moderate | Weak | Moderate |
| Security & Privacy | Required | Integrated | Integrated | Implied | Moderate |
| Vendor Governance | Required | Required | Expected | N/A | N/A |

## AI GOVERNANCE PROTOCOL BLUEPRINT MATRIX

*(Escalation, Intervention, Fallback, Traceability Requirements)*

| Protocol Requirement | What Must Be Defined | What IA Should Review | What Good Looks Like |
|---|---|---|---|
| Escalation Paths | Thresholds requiring review | Escalation SOPs | Clear triggers, roles, timelines |
| Human Intervention | When humans override AI | Override logs | Overrides documented + justified |
| Fallback Mechanisms | What happens when AI fails | Failover design | Safe defaults, manual alternative |
| Decision Logging | Tracking human + AI decisions | Log completeness | Timestamps, actor identification |
| Exception Handling | AI errors, anomalies, uncertainty | Exception workflow | Closed-loop resolution |
| Model Shutdown Triggers | Conditions to disable AI | Shutdown playbook | Accuracy drops, drift alerts |
| Mixed-Agent Controls | Who does what | RACI for human vs AI | No gaps or duplicate steps |

## AI GOVERNANCE MATURITY MATRIX (5 Levels)

*(Governance, Documentation, Controls, Monitoring, Oversight)*

| Level | Description | Indicators |
|---|---|---|
| 1. Ad Hoc | No formal AI governance | Shadow AI, undocumented models |
| 2. Emerging | Some guidelines, limited oversight | Basic policy, inconsistent monitoring |
| 3. Defined | Governance roles, processes documented | Model cards, data lineage, oversight |
| 4. Managed | Integrated lifecycle governance | Drift monitoring, retraining, KPIs |
| 5. Optimized | Enterprise AI governance system | ISO 42001 maturity, ethics board, audits |

## MIXED HUMAN–AI RESPONSIBILITY MATRIX

*(Who makes decisions, who reviews them, who overrides them)*

| Activity | AI Role | Human Role | IA Governance Expectation |
|---|---|---|---|
| Classification/Scoring | Generate predictions | Validate accuracy | Review accuracy KPIs |
| Decisions | Recommend | Approve/Decline | Confirm manual review |
| Exceptions | Flag anomalies | Investigate | Ensure workflow routing |
| Overrides | Identify uncertainty | Execute override | Logs audited |
| Monitoring | Detect drift | Review alerts | Trend analysis |
| Retraining | Trigger suggestion | Approve retraining | Governance approval |

## ASSURANCE EVIDENCE QUALITY MATRIX

*(What good vs weak evidence looks like)*

| Evidence Type | Strong Evidence | Weak Evidence |
|---|---|---|
| Data Governance | Lineage diagrams, DQ logs | Spreadsheet with raw fields |
| Model Development | Version control, model card | Screenshots, emails |
| Validation | Signed validation report | "We eyeballed it" |
| Oversight | Override logs, SOPs | Manual notes |
| Monitoring | Drift dashboards | Informal observations |
| Vendor Governance | Model cards, SOC2 | Marketing PDFs |

## AI LIFECYCLE GOVERNANCE MATRIX

*(Controls required at each lifecycle stage)*

| Lifecycle Stage | Required Controls | IA Expectations |
|---|---|---|
| *Design* | Requirements, ethics review | Validate traceability |
| *Development* | Versioning, documentation | Confirm reproducibility |
| *Testing* | Accuracy, fairness checks | Review test coverage |
| *Deployment* | Approval, change control | Confirm governance sign-off |
| *Monitoring* | Drift, KPIs, alerts | Reperform drift checks |
| *Retraining* | Trigger criteria | Verify model updates |
| *Retirement* | Decommissioning | Confirm data disposal |

**Disclaimer:**

The information provided in this training session and accompanying handouts is for educational purposes only. While every effort has been made to ensure the accuracy and completeness of the content, the presenter assumes no responsibility for errors, omissions, or any outcomes related to the application of the information provided. Participants are encouraged to seek professional advice or consult relevant guidelines for specific situations.