# Model AI Legislation Framework

# Tier 3 - Adoption & Implementation Guidance

Version 2.1 | December, 2025

**Practical pathways for integrating risk-based AI oversight into legislation, regulation, and governance**

Developed by AI Safety International

---

## 1.0 Purpose of Tier 3

Tier 3 provides practical guidance for adopting and implementing the Model AI Legislation Framework established in Tier 1 and supported by the technical methods described in Tier 2.

This tier is intended to assist legislators, legislative staff, regulatory agencies, and oversight bodies in understanding how a risk-based artificial intelligence oversight framework may be incorporated into real-world governance structures without excessive rigidity, technological lock-in, or government overreach.

Tier 3 does not propose statutory language, mandate specific regulatory designs, or prescribe technical controls. Instead, it outlines implementation pathways that preserve legislative authority while enabling consistent, enforceable, and adaptable oversight of artificial intelligence systems that pose meaningful risk to the public.

---

## 2.0 Intended Audience and Use

Tier 3 is written for:

• Legislators and legislative counsel
• Congressional and committee staff
• Regulatory agencies and rulemaking bodies
• Oversight authorities and review panels
• Policy analysts and auditors

This document is **non-binding** and **informative**. Its role is to reduce ambiguity during adoption and implementation by clarifying how the principles established in earlier tiers may be operationalized responsibly.

Nothing in this tier creates legal authority, regulatory obligations, or enforcement powers on its own. Any binding effect arises solely from enacted legislation or adopted regulation.

---

## 3.0 Mandatory Compliance and Flexible Implementation

The Model AI Legislation Framework is founded on a clear distinction:

**Compliance obligations are mandatory.**
**Implementation methods are flexible.**

Where Tier 1 establishes that covered systems must undergo risk-based assessment, documentation, and review, Tier 3 clarifies that legislatures and regulators may select from multiple legitimate approaches to achieving those requirements.

This structure ensures that safety expectations are enforceable, while preserving discretion in how compliance is demonstrated.

---

## 4.0 Legislative Adoption Models

Legislatures may adopt the framework through several viable pathways, depending on jurisdictional structure, regulatory capacity, and policy priorities.

### 4.1 Principle-Based Legislative Adoption

Under this model, legislation incorporates Tier 1 concepts directly by requiring that covered AI systems:

• undergo documented risk assessment prior to deployment
• identify and evaluate plausible failure modes
• implement proportional mitigation strategies
• retain documentation for review

Technical detail is intentionally excluded from statute and delegated to agencies, standards bodies, or rulemaking processes.

This approach promotes durability and avoids frequent statutory revision as technology evolves.

### 4.2 Reference-Based Legislative Adoption

Legislation may require compliance with recognized risk-assessment processes without naming specific methodologies.

For example, statutes may reference:

• structured risk assessment
• failure-mode analysis
• documented evaluation of severity, likelihood, and detectability

Agencies are then empowered to recognize acceptable technical standards or methodologies, including but not limited to AI-adapted FMEA models.

This approach preserves flexibility while maintaining enforceable obligations.

### 4.3 Phased or Transitional Adoption

Jurisdictions may introduce requirements through phased implementation, such as:

• initial grace periods
• pilot programs
• sector-specific rollout
• sunset or review clauses

Phased adoption may be appropriate where regulatory capacity is developing or where deployment contexts vary widely. Transitional flexibility applies to **timing**, not to the existence of the obligation itself.

---

## 5.0 Regulatory and Agency Implementation Pathways

Once legislative obligations are established, regulatory agencies may operationalize requirements through guidance, certification recognition, and review procedures.

### 5.1 Rulemaking and Interpretive Guidance

Agencies may issue guidance describing:

• documentation expectations
• review triggers
• reporting or retention standards
• criteria for demonstrating good-faith compliance

Guidance should focus on **process and outcomes**, not on content, internal model mechanics, or continuous monitoring of users.

## 5.2 Use of Technical Standards

Agencies may recognize external technical standards or methodologies developed by professional or international standards bodies.

This approach allows technical requirements to evolve without statutory revision and reduces the risk of regulatory stagnation.

Recognition of standards does not require exclusivity; equivalent methodologies may be accepted where they meet comparable analytical rigor.

## 5.3 Oversight Without Surveillance

Implementation under this framework does not require:

• monitoring private user conversations
• inspecting internal model architectures
• reviewing training data or proprietary systems

Oversight relies on documented assessments, reviewable processes, and incident-based evaluation rather than continuous surveillance.

---

# 6.0 Role of Certification in Compliance

Certification serves as the primary mechanism for demonstrating compliance with Tier 1 obligations.

Certification within this framework:

• evaluates whether required processes were followed
• verifies documentation and review practices
• supports accountability without prescribing design choices

Certification is not an approval of ideas, content, or intent. It is an assessment of risk management discipline.

## 6.1 Certification Weight and Effect

Certification carries substantial practical weight by:

- enabling market access
- supporting procurement eligibility
- informing insurer risk evaluation
- serving as evidence of due diligence

Failure to certify where required may result in regulatory consequences, market exclusion, or increased liability exposure.

## 6.2 Distributed Certification Authority

Certification is not centralized in a single governmental body.

Instead, responsibility is distributed among:

- recognized certification organizations
- independent auditors
- standards bodies
- oversight agencies

Government's role is to require certification and recognize acceptable certifying processes, not to perform certification directly.

---

# 7.0 Use of AI-FMEA and Equivalent Methods

AI-FMEA is presented as one recognized structured methodology for conducting risk-based assessment.

Its role within implementation is:

- illustrative, not exclusive
- acceptable, not mandatory
- supportive, not determinative

Organizations may use AI-FMEA or alternative methodologies that demonstrate comparable rigor in identifying, evaluating, and mitigating risk.

Regulators and courts may consider documented assessments as evidence of good-faith compliance without requiring adoption of any specific tool.

---

## 8.0 Compliance, Review, and Enforcement Alignment

### 8.1 Proportional Compliance Expectations

Compliance expectations scale with risk.

Systems exhibiting greater relational influence, continuity, personalization, or user reliance warrant deeper assessment and stronger safeguards. Lower-risk systems may require lighter documentation and periodic review.

### 8.2 Incident Response and Accountability

When harm occurs, documented assessments provide a basis for determining whether risks were:

• foreseeable
• identified
• reasonably mitigated

Distinguishing between unforeseeable failure and negligent omission supports fair enforcement and discourages reckless deployment.

### 8.3 Safe Harbor for Good-Faith Compliance

Good-faith adherence to assessment and documentation requirements may provide mitigation of liability or enforcement severity.

This incentivizes responsible behavior without shielding actors from accountability for willful disregard or falsification.

---

## 9.0 Cross-Jurisdiction and International Compatibility

The framework is designed to coexist with:

• federal and state regulatory systems
• sector-specific oversight regimes
• international risk-based standards

Market access and certification recognition provide alignment incentives without extraterritorial enforcement.

## 10.0 Maintaining Innovation and Investment Stability

Risk-based oversight supports innovation by:

• providing predictable expectations
• reducing regulatory uncertainty
• avoiding technology-specific mandates
• protecting public trust

Innovation remains unrestricted outside defined risk thresholds. Safety requirements apply where harm potential is material and foreseeable.

---

## 11.0 Common Implementation Pitfalls (Advisory)

Experience from other safety-critical domains suggests avoiding:

• embedding technical detail in statute
• mandating single tools or vendors
• equating oversight with surveillance
• ignoring cumulative or relational effects
• treating documentation as a formality

Tier separation helps prevent these failures.

---

## 12.0 Preparing for Evolution

Artificial intelligence systems, deployment contexts, and risk profiles will continue to evolve.

Periodic review of implementation practices, recognition of updated standards, and reassessment of scope thresholds are essential to maintaining relevance without destabilizing governance structures.

---

## 13.0 Closing Note

Tier 3 completes the Model AI Legislation Framework by translating principles and technical foundations into practical governance pathways.

Together:

• Tier 1 defines **what must be addressed**
• Tier 2 explains **how risks may be evaluated**
• Tier 3 guides **how oversight may be implemented**

This framework does not seek to restrain progress, but to ensure that systems capable of influencing human judgment, behavior, and wellbeing are deployed with responsibility equal to their power.

---

**End of Tier 3 — Adoption & Implementation Guidance**