

Artificial Intelligence Governance Framework in Privacy

Developed by: **Fernanda Sabatine**



Professional with strategic expertise in data protection and information security. Created the 4+1 Governance Framework in AI as a tool to integrate privacy into strategic decisions involving AI models, combining regulations such as GDPR, LGPD, and relevant international references.

Concept: 4+1 Governance Framework

Objective: Provide a practical and adaptable governance structure to help organizations develop, implement, and manage Artificial Intelligence systems ethically, securely, and in compliance with key data protection regulations (LGPD, GDPR) and emerging AI legislation (especially the EU AI Act).

PILLAR 1: GOVERNANCE AND ACCOUNTABILITY

(The Strategic Foundation)

Strategic Objective: Establish a clear framework of responsibility and oversight for all AI initiatives, ensuring top management accountability and conscious, well-documented decision-making. This pillar aligns with ISO 42001 for an AI Management System.

Essential Controls (Practical Actions):

- **Definition of Roles and Responsibilities:** Formally designate those responsible for AI governance, including a multidisciplinary AI Ethics Committee and project leaders for each AI system.
- **Artificial Intelligence Policy:** Develop and approve a high-level policy outlining the organization's ethical principles (fairness, transparency, security), permitted and prohibited use cases, and approval procedures for new AI initiatives.
- **Inventory of AI Systems (AI ROPA):** Maintain a detailed registry of all AI systems in use or under development, documenting their purpose, training and operational data, decision logic, and associated risks (an extension of the privacy ROPA) [*The Intersection of AI and Governance – IAPP, p. 10*].
- **Risk Analysis (EU AI Act):** Classify each AI system according to risk levels defined in the EU AI Act (Unacceptable Risk, High Risk, Limited Risk, Minimal Risk) to determine applicable compliance obligations [*EU AI Act – Compliance Matrix – IAPP.pdf, p. 4, AI ACT*].

PILLAR 2: DATA AND MODEL LIFECYCLE MANAGEMENT

(From Conception to Disposal)

Strategic Objective: Ensure that privacy and security principles are applied throughout the lifecycle of an AI system, from data collection for training to model deactivation, incorporating Privacy by Design.

Essential Controls:

- **Training Data Governance:**
 - *Legality and Quality:* Ensure that data used for training is lawfully obtained, with the correct legal basis (LGPD/GDPR), and is high quality, representative, and relevant to the model's purpose [*Legal Basis for Data Protection in the Use of AI – Baden Württemberg, p. 5*].
 - *Data Minimization:* Use only data strictly necessary for training, applying anonymization or pseudonymization techniques wherever possible.
- **Impact Assessment (AIA / DPIA for AI):** Conduct an Algorithmic Impact Assessment (AIA), an expanded version of RIPD/DPIA. It should assess privacy risks, model bias, discrimination, fairness, and social/ethical impacts [*AI and Algorithms in Risk Assessment – ELA, p. 14*].
- **Privacy by Design in AI Engineering:** Incorporate privacy and security requirements from the design phase, making data protection a native component.
- **Model Lifecycle Management:** Establish processes for continuous monitoring of model performance, periodic recalibration, and secure deactivation at end-of-life.

PILLAR 3: FAIRNESS, TRANSPARENCY, AND EXPLAINABILITY (XAI)

(The Pillar of Trust and Ethics)

Strategic Objective: Address algorithmic bias, ensure automated decisions can

be explained to affected individuals, and be transparent about how and when AI systems are used — in line with ICO and other authorities [*Regulating AI – The ICO’s Strategic Approach*].

Essential Controls:

- **Bias Mitigation (Fairness):** Apply technical tests during development to detect and mitigate bias in training data and algorithm behavior. Document these tests to demonstrate due diligence [*Generative AI Controls Framework – IBM, p. 12*].
- **Explainability Mechanisms (XAI):** For systems making significant decisions about individuals, implement tools that explain the logic behind each decision in plain language [*Artificial Intelligence and Privacy – Solove, p. 19*].
- **Meaningful Human Oversight:** Ensure that high-risk decisions are never fully automated. There must be clear processes for human review and intervention (human-in-the-loop), particularly for legally impactful decisions.
- **Transparency of Use:** Clearly inform individuals when interacting with AI systems (e.g., chatbots) or when decisions affecting them are based on automated processing — as required by the EU AI Act for limited-risk systems.

PILLAR 4: AI SECURITY AND RESILIENCE

(Protection Against Emerging Threats)

Strategic Objective: Protect AI systems against ecosystem-specific threats that go beyond traditional information security, ensuring the integrity, confidentiality, and availability of models and data.

Essential Controls:

- **Training Data Protection:** Apply robust security controls to protect training datasets from corruption or unauthorized access, preventing data poisoning.
- **Defense Against Adversarial Attacks:** Implement techniques to make AI models more resilient against adversarial inputs designed to mislead the system [*Agentic AI – Threats and Mitigations – OWASP*, p. 6].
- **AI Infrastructure Security:** Secure the entire AI pipeline — from development and training to production environments (APIs, servers), including vulnerability management for open-source libraries.
- **AI Incident Response Plan:** Adapt the organization's incident response plan to include AI-specific scenarios, such as adversarial attacks, mass-biased outcomes, or model leakage.

PILLAR 5: PRIVACY AS A STRUCTURAL AXIS OF AI GOVERNANCE

(The Ethical and Legal Foundation of Automated Systems)

Data protection is not just a technical or regulatory measure: it is a fundamental

component of responsible AI systems, ensuring respect for human dignity, user autonomy, and social transparency.

Strategic Objective:

Incorporate personal data protection principles and requirements across the lifecycle of AI systems — from design to deployment and deactivation — ensuring compliance with LGPD, GDPR, and EU AI Act when applicable.

Guiding Principles:

- Privacy by Design
- Data Minimization
- Specific and Legitimate Purpose
- Informed and Freely Given Consent
- Priority to Data Subject Rights
- Transparency and Accountability

Essential Controls:

- **Integration with Privacy Team and DPO:** Ensure privacy experts are actively involved in AI projects from ideation onward.
- **PIA and DPIA Integrated with AI Governance:** Conduct Privacy Impact Assessments alongside EU AI Act risk classification, focusing on automated decisions, sensitive data, and risks to individual freedoms.
- **Extended AI ROPA with Privacy Focus:** Include fields such as legal basis, data types, anonymization/pseudonymization processes, and data recipients.

- **Transparent and Auditable Consent Management:** Implement mechanisms for consent collection and management in AI systems using personal data, with full traceability.
- **Fair and Non-Discriminatory Processing:** Adopt clear anti-discrimination policies, validating models from privacy and equality perspectives.
- **Automated Data Subject Rights:** Develop operational methods to ensure:
 - Explanation of automated decisions
 - Human review of critical decisions
 - Portability of data used in AI
 - On-demand deletion or anonymization
- **Organizational Training and Privacy Culture in AI:** Launch regular training programs for developers, leaders, and ethics committees focused on applied privacy in AI.
- **Implementation of Privacy Enhancing Technologies (PETs):** Use techniques like differential privacy, federated learning, and homomorphic encryption to train and operate models without exposing raw personal data.
- **Generative AI Governance:**
 - *Acceptable Use Policy:* Create clear guidelines for employees regarding what types of information (especially personal or confidential data) may or may not be entered into third-party generative AIs (e.g., ChatGPT, Gemini, Copilot, DeepSeek).

- *Hallucination Analysis*: Establish procedures to verify the accuracy of AI-generated information — especially when involving personal data — to prevent storing or creating false outputs.
- **Data Subject Rights Management in AI Systems**: Define clear procedures to fulfill access, rectification, and especially deletion rights, including the complex task of removing a specific data's influence from an already trained model (algorithmic right to be forgotten).
- **Legal Basis Adequacy**: Conduct a rigorous analysis to ensure the chosen legal basis (e.g., consent, legitimate interest) is appropriate and defensible for each AI system's specific purpose, avoiding common GDPR/LGPD compliance pitfalls.