# Generative AI Governance Implementation Guide

**Type:** Implementation Guide

**Target Audience:** CIOs, AI Program Managers, Governance Leaders

**Timeline:** 6-12 Months

## 1. Executive Summary

### The Generative AI Challenge

Generative AI (GenAI) differs fundamentally from traditional predictive AI. While traditional AI classifies existing data (e.g., "is this transaction fraudulent?"), GenAI creates new content (text, code, images). This shift introduces novel risks that legacy governance frameworks often miss, such as **hallucinations** (confident falsehoods), **copyright infringement** (in training data and outputs), and **shadow proliferation** (employees using consumer tools for sensitive enterprise data).

### Strategic Objectives

This guide provides a roadmap to move your organization from ad-hoc, risky adoption to a managed, mature state. Implementing this framework aims to achieve:

1. **Risk Reduction:** Mitigating IP leakage and reputational damage from unverified outputs.

2. **Regulatory Readiness:** Preparing for obligations under the EU AI Act (GPAI provisions) and emerging US state laws.

3. **Operational Efficiency:** Reducing "Shadow AI" usage by providing safe, sanctioned alternatives.

### Implementation Timeline Overview

| Phase | Focus | Duration |
|---|---|---|
| **1. Foundation** | Discovery, Ownership, & Immediate Guardrails | Weeks 1-4 |
| **2. Policy** | Use Case Definitions & Data Standards | Weeks 5-8 |
| **3. Controls** | Technical Gateways & Monitoring | Weeks 9-12 |
| **4. Culture** | Training, Literacy & Champions | Weeks 13-16 |
| **5. Operations** | Monitoring, Auditing & Iteration | Ongoing |

# Phase 1: Foundation (Weeks 1-4)

*Goal:* *Stop the bleeding on data leakage and establish clear accountability.*

## 1.1 "Shadow AI" Discovery Audit

You cannot govern what you do not know. Most organizations have significantly more GenAI usage than IT realizes.

- **Survey Business Units:** Ask specifically about use of tools like ChatGPT, Claude, Jasper, or GitHub Copilot.
- **Review Vendor Contracts:** Check existing SaaS platforms (CRM, HRIS) for "AI features" enabled by default.
- **Network Analysis:** Work with security to identify traffic to common AI domains.

## 1.2 Establish Governance Committee

Do not create a siloed "AI Department." Establish a cross-functional steering committee.

- **Chair:** Executive Sponsor (CIO/CDO/CRO).
- **Core Members:** Legal (IP/Copyright), CISO (Data Security), HR (Employee Policy), and Data Privacy Officer.
- **Cadence:** Weekly during implementation; monthly thereafter.

## 1.3 Quick Wins Checklist

> **Strategic Insight:** Don't wait for a perfect policy to issue guidance. Issue an "Interim Guidance" memo immediately.

- ☐ Issue Interim Guidance: Explicitly state which public tools are prohibited for internal data.
- ☐ Block High-Risk Domains: If no enterprise agreement exists, block traffic to consumer AI sites.
- ☐ Identify "Red" Data: List specific data types (PII, trade secrets, source code) that never go into a public prompt.

**Phase 1 Deliverables**

- ☐ Inventory of current GenAI usage (Authorized vs. Shadow).
- ☐ AI Steering Committee Charter ratified.
- ☐ Interim Acceptable Use Memo distributed to all staff.

# Phase 2: Policy Development (Weeks 5-8)

**Goal:** *Define the "rules of the road" for acceptable use and procurement.*

## 2.1 Use Case Classification Framework

Adopt a risk-tiered approach to avoid stifling low-risk innovation.

| Risk Tier | Definition | Examples | Governance Requirement |
|---|---|---|---|
| **Prohibited** | Unacceptable risk to rights, safety, or IP. | Generating deepfakes without consent; automated hiring decisions without human review. | **Not Allowed** |
| **High** | Customer-facing or high-stakes operational impact. | Medical advice; legal contract drafting; automated code generation for production. | **Strict:** Legal review, red-teaming, mandatory human-in-the-loop. |
| **Medium** | Internal operational efficiency. | Summarizing internal meeting notes; drafting marketing copy; internal knowledge search. | **Moderate:** Human review required; output verification. |
| **Low** | Individual productivity; no sensitive data. | Brainstorming ideas; Excel formula generation; drafting generic emails. | **Light:** User training; acceptable use policy adherence. |

## 2.2 Data Handling Requirements

Define how data interacts with models based on the AI Technology Stack.

- **Public/Consumer Models:** No proprietary data, PII, or IP allowed.
- **Enterprise/Private Instances:** Proprietary data allowed *if* the vendor contract guarantees zero data retention for training base models (e.g., Azure OpenAI, AWS Bedrock).
- **Self-Hosted/Open Source:** Highest sensitivity data allowed; requires security audit of the hosting infrastructure.

## 2.3 Vendor Evaluation Criteria

Update procurement questionnaires to include GenAI-specific checks:

- Does the vendor use customer data to train their foundation models? (Must be "No" for sensitive data).
- Can the vendor indemnify the organization against IP infringement claims on outputs?
- Is the model statically versioned, or does it update automatically (causing drift)?

**Phase 2 Deliverables**

☐ GenAI Risk Classification Matrix.
☐ Updated Data Classification Policy to include "AI Prompts."
☐ Procurement addendum for AI vendors.

# Phase 3: Technical Controls (Weeks 9-12)

*Goal: Enforce policy through technology, not just trust.*

## 3.1 The AI Gateway Strategy

Implement an API Gateway or "Wrapper" application for employee access to LLMs. This prevents direct API key sharing and centralizes logging.

- **Centralized Access:** Employees log in via SSO to a corporate portal (e.g., "CorporateGPT") rather than using public accounts.
- **Cost Controls:** Set token limits per user/department to prevent budget overruns.

## 3.2 Input and Output Guardrails

- **PII Filtering (Input):** Implement DLP (Data Loss Prevention) scanners on the prompt input field to block credit cards, SSNs, or specific project codenames.
- **Toxic Content Filtering (Output):** Configure content filters (e.g., Azure Content Safety) to block hate speech, violence, or jailbreak attempts.
- **Prompt Logging:** Log all prompts and responses for audit purposes (ensure employees know this is not private).

## 3.3 Integration Standards

Define how GenAI connects to other systems.

- **Sandboxing:** GenAI code interpreters must run in isolated containers, not on the user's local machine.
- **Human-in-the-Loop (HITL) UI:** Interfaces must explicitly label AI-generated content and require a human "verify and approve" step before content is published or code is merged.

**Phase 3 Deliverables**

☐ Corporate GenAI Portal/Gateway deployed.
☐ DLP rules updated for prompt injection and PII.
☐ Centralized logging dashboard for AI usage established.

# Phase 4: Training & Culture (Weeks 13-16)

*Goal: Mitigate the "Human Factor" risks (hallucinations and over-reliance).*

## 4.1 Role-Based Training Curriculum

Training cannot be one-size-fits-all. Segment your audience:

### Tier 1: All Employees (The Basics)

- What is Generative AI? (It predicts words, it doesn't "know" facts).
- The risks of hallucination (Verify everything).
- Data privacy: "If you wouldn't post it on Reddit, don't put it in a public LLM."

### Tier 2: Managers & Content Creators

- How to review AI work.
- Copyright implications of AI-generated assets.
- Prompt engineering basics for efficiency.

### Tier 3: Developers & Data Scientists

- Secure coding with AI assistants.
- Preventing prompt injection.
- Model evaluation metrics.

## 4.2 "GenAI Champions" Program

Identify power users in various departments (Marketing, HR, Dev).

- **Role:** Act as the first line of defense and innovation.
- **Responsibility:** Vet use cases in their department before bringing them to the Governance Committee.
- **Benefit:** They get early access to new tools and features.

### Phase 4 Deliverables

☐ Training modules launched in LMS.
☐ "GenAI Champions" identified and onboarded.
☐ Internal communications campaign regarding responsible use.

# Phase 5: Monitoring & Iteration (Ongoing)

*Goal: Ensure the framework evolves with the technology.*

## 5.1 Key Performance Indicators (KPIs)

- **Adoption Rate:** % of employees using sanctioned tools vs. estimated shadow usage.
- **Incident Rate:** Number of data leakage incidents or hallucination-related errors reported.
- **Review Velocity:** Average time to approve a new GenAI use case (Target: <2 weeks).

## 5.2 Incident Tracking

Treat AI incidents distinct from standard IT incidents. Create specific categories for:

- **Hallucinations:** AI inventing facts in business contexts.
- **Bias:** AI generating discriminatory output.
- **Prompt Injection:** External attempts to manipulate public-facing bots.

## 5.3 Regulatory Watch

Assign a legal/compliance owner to monitor:

- **EU AI Act:** Watch for "General Purpose AI" obligations.
- **Copyright Law:** Monitor court cases (e.g., NYT v. OpenAI) for shifts in IP liability.

### Phase 5 Deliverables

☐ Quarterly AI Governance Report to the Board.
☐ Bi-annual review of the "High Risk" use case list.

# Appendices

## Appendix A: RACI Matrix for GenAI Governance

| Activity | AI Steering Comm. | IT / CISO | Legal / Privacy | Business Unit Lead | End User |
|---|---|---|---|---|---|
| Define Acceptable Use Policy | A | C | R | C | I |
| Approve High-Risk Use Case | A | C | C | R | I |
| Vendor Security Assessment | I | A | C | I | - |
| Prompt Engineering | - | - | - | A | R |
| Verify AI Output Accuracy | - | - | - | A | R |
| Report AI Incident | I | C | C | C | R |

## Appendix B: Sample Acceptable Use Policy (Excerpt)

***Principle of Human Accountability:***
*"AI systems are tools to assist, not replace, human judgment. You are responsible for the accuracy, legality, and appropriateness of any work you produce, regardless of whether you used AI to draft it.* ***You must review and verify all AI outputs.****"*

***Principle of Data Confidentiality:***
*"Do not input Restricted or Confidential data into any AI tool that has not been explicitly approved by Information Security. Assume all inputs to public chatbots (e.g., free ChatGPT) become public information."*

## Appendix C: GenAI Risk Assessment Matrix

*Use this to score new proposals.*

| Risk Impact Factors | Low (1) | Medium (2) | High (3) |
|---|---|---|---|
| **Data Sensitivity** | Public data only | Internal-only data | PII, PHI, or IP |
| **Human Oversight** | Output verified by human | Human reviews samples | Fully automated action |
| **Failure Consequence** | Minor inconvenience | Reputational risk | Legal liability / Safety |
| **Total Score** | 3-4 (Green) | 5-7 (Yellow) | 8-9 (Red) |

- **Green (3-4):** Proceed with standard guidelines.
- **Yellow (5-7):** Requires Risk Assessment & Manager Approval.
- **Red (8-9):** Requires Steering Committee Approval & Red Teaming.