

Day 3



Generative AI Security

Risk-Based Approach to AI Governance:

What Is an AI Risk?

An AI risk is any possibility that an AI system may cause harm, fail, or create negative impact for users, business, or society. Simply, it talks about “*What bad thing can happen because of this AI system?*”

Categories of AI risks: (Generated using Gemini)

Categories of AI Risks (Very Easy Breakdown)

There are six major categories of risks in AI systems.

1. Data Risks Biased data, Incorrect or low-quality data, Sensitive data leakage, Data poisoning attacks, Unauthorized access, Lack proper consent. Impact: results become unfair, wrong, or insecure.	2. Model Risks Model bias, Hallucinations, Lack explainability, Overfitting / underfitting, Unreliable predictions, Unreliable predictions, Poor performance in real-world scenarios. Impact: affects real-world scenarios.	3. Operational Risks Misconfigured pipelines, Incorrect model updates, No monitoring or alerts, Human error, Integration failures, Incomplete documentation, No fallback backups, Outages, Impact: system breakdowns, unpredictable behavior.
4. Security Risks Data poisoning, Prompt injection, Adversarial examples, No monitoring examples, Model theft (IP leak), Jailbreaking LLMs, Unauthorized agent actions. Impact: attacker controls the AI or extracts sensitive information.	5. Ethical Risks Discrimination, Harmful recommendations, Manipulative content, Invasion of privacy, Unsafe autonomy or privacy, Negative societal impact. Impact: unethical or unsafe behavior towards individuals and society.	6. Regulatory Risks Violating the EU AI Act, Breaking privacy laws (GDPR / DPPA Act), Ignoring ISO 27001 requirements, Not having documentation for audits without safeguards, legal penalties, fines, and reputation damage.

Risk Assessment Techniques:

AI risk assessment evaluates how dangerous an AI system could be. Risk assessment can be qualitative (Uses words, not numbers), quantitative (Uses numbers, scoring, metrics, and probabilities), or scoring-based.

Risk Prioritization & Heat Maps:

Heat maps help prioritize which risks need immediate attention. After scoring risks, we organize them in a heat map:

Likelihood	Impact	Risk Level
High + High	🔥 Red	Critical
High + Medium	🟠 Orange	High
Medium + Medium	🟡 Yellow	Medium
Low + Low	🟢 Green	Low

Simply,

Red risks → fix immediately

Green risks → normal monitoring

Yellow → moderate controls

Linking Risks to Controls and Mitigations:

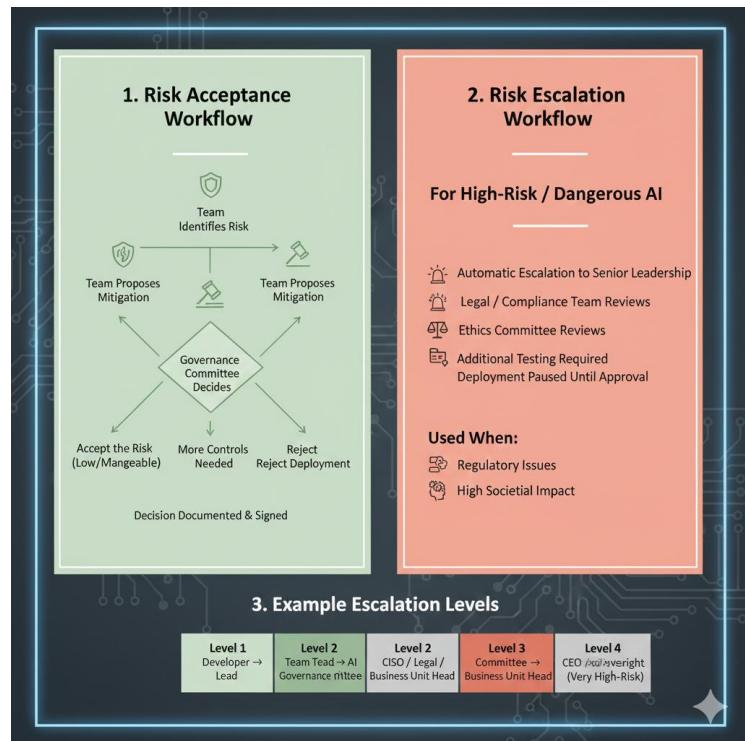
For every risk identified, governance requires specific controls (A control is a safeguard, procedure, or mechanism put in place to reduce, prevent, or manage a specific risk.). Governance ensures every risk has a matching control.

Example Mapping:

Risk	Control / Mitigation
Biased data	Fairness testing, rebalancing datasets
Model drift	Continuous monitoring, retraining
Prompt injection	Strong input validation, guardrails
Model theft	API rate limits, access control, encryption
Privacy leakage	Differential privacy, anonymization
Non-compliance	Documentation, human oversight, audit logs

Governance Workflows for Risk Acceptance & Escalation:

If a risk cannot be fully removed, governance defines what to do. Governance workflows define who can accept, reject, or escalate high-risk AI decisions.



AI Stakeholders & Their Roles:

AI lifecycle and stakeholders:



- Many people work together to make AI safe.
- ML engineers build models, data teams prepare data, security protects systems.
- Compliance ensures laws are followed, business owners justify use cases.
- Executives give final approval, users provide feedback.

RACI Models for AI Projects:

RACI = Responsible, Accountable, Consulted, Informed

Typical matrix:

Activity	Responsible	Accountable	Consulted	Informed
Data collection	Data team	Data owner	Legal/Compliance	Executives
Model training	ML Engineers	AI Lead	Security, Data	Business owner
Security testing	Security team	CISO	ML Engineers	Executives
Risk assessment	Compliance	CISO/CTO	ML + Security	Business owner
Deployment	ML Ops	CTO	Security + Compliance	All stakeholders
Monitoring	Security + ML Ops	CISO	Business	Regulators (if needed)

RACI helps avoid confusion about who does what.

Note: Collaboration happens during training, deployment, and monitoring.

--The End--