# Data Security within an AI Environment

# Data Security within AI Environments

## Working Draft

## Authors:

- Login to your Google account to access full editing permission.
- Change from Editing to Suggesting in the upper right of the Google doc for tracking each author's edits.

*Please contact research-support@cloudsecurityalliance.org to request full access to author this document.*

## Reviewers/Visitors:

- If you have a Google Account, please login before commenting. Otherwise, please note your name and affiliation in the comment you leave.
- Use the Comments or Suggesting features on Google docs to leave your feedback on the document. Suggestions will be written in and identified by your Google Account. To use the comments feature, highlight the phrase you would like to comment on, right click and select "Comment" (or Ctrl+Alt+M). Or, highlight the phrase, select "Insert" from the top menu, and select "Comment." All suggestions and comments will be reviewed by the editing committee.

*For more information about Google's Comments feature, please refer to http://support.google.com/docs/bin/answer.py?hl=en&answer=1216772&ctx=cb&src=cb&cbid=-rx63b0fx4x0v&cbrank=1*

The permanent and official location for the Data Security Working Group is
https://cloudsecurityalliance.org/research/working-groups/data-security

# Acknowledgments

## Lead Authors

Rocco Alfonzetti
Alex Kaluza
Vashti Horvát
Oliver Forbes
Gopi Ramamoorthy
Onyeka Illoh
Ikechukwu Okoli
Prateek Mittal
Mahesh Adulla

## Contributors

Naveed Afzal
Andres Castagna
Jan Gerst
Swaminadhan Jagadeesan
Richard Kabanda
Sachindra Narayan
Paul Son
Saurabbh Srivastava
Yuanji Sun
Sarah Templey
Ayoob Ullah
Tuleun Washima

## Reviewers

Mahesh Adulla
Kachi Agu
Jayesh Dalmet
Rob Doyon
Udo Duro
Akshatha Gangadharaiah
Hariprasad Holla
Will Lemos
Ramesha Reddy Thimmasandra
Nsikak-Abasi Shammah Una
Akshat Vashishtha

Chad Walter
Washima Tuleun

## CSA Global Staff

Hillary Baron
Alex Kaluza

# Table of Contents

# Introduction

This paper explores the evolving landscape of data security in artificial intelligence (AI) environments and provides practical guidance aligned with the Cloud Security Alliance (CSA) AI Controls Matrix (AICM). As AI systems increasingly depend on large language models (LLMs), decentralized learning architectures, and privacy-enhancing technologies (PETs), traditional data protection practices must adapt to address new threats and operational realities.

The foundational principles of data security—confidentiality, integrity, and availability—remain essential, but they must be applied differently in modern AI systems. Distributed processing, challenges confidentiality, since sensitive data may be exposed across multiple environments. Integrity can be compromised by adversarial inputs and manipulated training data. Availability must be preserved while managing potential misuse or unintended system behavior. Addressing these risks requires AI-specific approaches to data protection.

To support this effort, this paper maps existing AICM controls (DSP-01 through DSP-24) to the AI data lifecycle and identifies areas where current safeguards are insufficient. It proposes four new controls for future inclusion in the AICM: Prompt Injection Defense, Model Inversion and Membership Inference Protection, Federated Learning Governance, and Shadow AI Detection. These additions are intended to fill critical gaps in the current framework and reflect emerging risks introduced by generative AI.

The emergence of multi-modal AI systems—which process text, images, audio, and video simultaneously—introduces unprecedented cross-modal data leakage risks. Information from one modality can inadvertently expose sensitive data from another, creating attack vectors that traditional security controls cannot address. Furthermore, the rise of agentic AI systems capable of autonomous decision-making and cross-boundary data access demands new containment strategies that extend beyond conventional perimeter security.

The CSA AI Controls Matrix provides a comprehensive framework for addressing these challenges through systematic risk management. Our approach incorporates industry best practices—including alignment with the NIST AI Risk Management Framework's core functions—while maintaining focus on cloud-specific implementations and practitioner needs. This ensures organizations can achieve both regulatory compliance and operational security through a single, integrated control set.

This guidance is intended for security architects, compliance leaders, and AI system owners who are working to operationalize data security and risk management in AI environments. It aligns with broader CSA resources, the OWASP Top 10 for LLMs, and regulatory requirements such as General Data Protection Regulation (GDPR) and Health Insurance Portability and Accountability Act (HIPAA). By addressing these challenges directly, this paper supports the development of trustworthy, secure, and resilient AI systems.

# AI Security Strategy Foundations

A data-centric AI security strategy prioritizes the protection of data across the entire AI lifecycle, beginning with collection and preprocessing and continuing through model training, inference, and storage. Because AI systems rely on the data they process for functionality and value, it is essential to ensure that this data remains confidential, accurate, and available throughout the lifecycle. Without these safeguards, the risk of unauthorized access, manipulation, and information leakage increases significantly.

Traditional AI security approaches focus on securing models, systems, and infrastructure. These methods defend against threats, such as adversarial attacks, data poisoning, model extraction, and unauthorized inference. Controls are typically applied at the network, application, and system layers to protect the broader environment where AI systems operate.

As AI deployments become more complex and distributed, many organizations are adopting a hybrid security strategy that combines traditional protections with data-centric safeguards. This approach offers end-to-end security across diverse environments, including on-premises data centers, cloud platforms, and edge devices. By addressing vulnerabilities throughout the full AI lifecycle, this integrated model helps ensure that AI operations remain secure, resilient, and compliant with evolving data protection requirements.

# Understanding AI and Its Intrinsic Demand for Data

Artificial intelligence depends heavily on data, which presents both opportunities and challenges for cybersecurity professionals. AI is transforming industries by enabling automation, improving decision-making, and delivering insights through large-scale data analysis. However, the need for massive, diverse data sets also introduces significant cybersecurity risks. Machine learning systems require high volumes of data for training, validation, and testing, and the quality and sensitivity of that data directly influence their reliability and security. As organizations increasingly adopt AI technologies, it is essential to understand the implications of data collection, storage, and processing from a security perspective.

# Data Collection

AI systems ingest data from a wide variety of sources. These include:

- **Sensors and IoT devices:** Collect real-time information from physical environments.
- **Mobile and Edge devices**: User generated behavioral and contextual data from smartphones, tablets, etc.
- **Social media platforms:** User-generated content and behavioral data are gathered.
- **Enterprise systems:** Provide transactional, financial, and operational data.
- **Public and open data sets:** Include structured and unstructured information for broad applications (e.g., census data, benchmarks, scientific data sets).
- **Cameras and video feeds**: Visual data for perception and recognition tasks (e.g., surveillance cameras, autonomous vehicles).
- **User Interactions**: Gather information through user interactions by  processing user input and respond appropriately.
- **Research Libraries**: Collaborative data repositories used for everything from market studies through to legal, medical, economic, and financial research.
- **Marketing and Sales Content**: Websites, brand guides, sales training, white papers, technical briefs, blogs, podcasts, product descriptions, images, and anything posted on and accessible through internet queries.
- **Synthetic data** generated from AI models.

These sources enrich training data but expand the attack surface.The origin, integrity, consent, and lawful basis associated with data collected from these sources must be carefully managed to avoid introducing security and privacy risks. Data collection must follow data minimization and purpose limitation principles, ensuring only necessary data is collected strictly for clearly defined, specific purposes.

## Data Annotation as a Security-Critical Activity

Data annotation represents a primary attack vector often overlooked in traditional security frameworks. This process, which transforms raw data into training-ready data sets, can be exploited through the following.

**Semantic Supply Chain Attacks**:
- Malicious annotators introducing systematic biases or backdoors through labeling
- Poisoning attacks where incorrect labels degrade model integrity before training begins
- Privacy breaches when annotators access un-redacted sensitive information

**Secure Annotation Controls**:

- Establish cryptographically signed data provenance tracking using blockchain or merkle trees, scalability trade-offs may exist (e.g., blockchain storage costs, annotation throughput) and alternatives like tamper-evident logs (e.g., AWS QLDB, Sigstore) may fill this gap.
- Implement inter-annotator agreement metrics (Cohen's kappa > 0.8 for critical domains)
- Deploy automated PII detection and redaction before human annotation, using tools like Microsoft Presidio or Google DLP API
- Require security clearances for annotators handling sensitive domains (medical, financial, defense)
- Create isolated annotation environments with data masking and audit logging

**Quality Assurance Framework**:
- Sample 10% of annotations for expert review
- Track annotator performance metrics and flag statistical anomalies
- Implement version control for annotation guidelines with change justification
- Maintain immutable audit logs of all annotation activities

# Data Storage

Once collected, data is typically stored in centralized systems, such as databases, data lakes, or cloud storage platforms. These repositories must support high-volume and high-velocity data flows while maintaining availability for AI processing tasks. However, their scale and accessibility also make them attractive targets for attackers. Secure storage design must include encryption, access controls, and data segregation based on sensitivity levels.

Associated risks in Data Storage include:

- Model inversion from stored data – Attackers reconstruct training data from models
- Unauthorized access – Misconfigured cloud storage or weak identity and access management (IAM)
- Insufficient encryption – Data at rest exposed to insiders or attackers
- Data residency issues – Storage in jurisdictions with conflicting regulations
- Shadow datasets – Forgotten or duplicate datasets becoming entry points
- Shadow AI / unsanctioned model use – Hidden storage and use of unapproved datasets or models

# Data Processing

The processing phase involves preparing data for use in AI models and includes several resource-intensive steps:

- **Data cleaning and preprocessing:** Ensures quality, consistency, and removal of errors or outliers.
- **Feature extraction:** Identifies the most relevant attributes for training models.
- **Model training and evaluation:** Uses the prepared data to develop and validate AI performance.

Each of these stages presents a potential point of exposure if controls are not in place. Temporary data artifacts, derived features, and intermediate outputs can all carry sensitive information and must be protected accordingly. This can be achieved by protecting temporary artifacts using encrypted scratch storage, access isolation, dataset versioning, and time-bound deletion (TTL).

# Cybersecurity Challenges

The volume and centralization of data in AI systems make them high-value targets for cyberattacks. A successful breach could expose personal, financial, or proprietary data, resulting in regulatory violations, reputational harm, or business disruption. AI systems are also vulnerable to adversarial attacks, where inputs are intentionally manipulated to trick models into producing incorrect outputs. Additionally, data poisoning attacks can compromise model integrity during training, and model theft can reveal intellectual property and business logic.

Users/Employees connecting to unapproved AI platforms and services create unmonitored risks. These risks are amplified by the widespread use of third-party models and open-source components.

# Mitigation Strategies

To address these challenges, organizations should implement a layered data security strategy tailored to AI workflows. This includes:

- Encryption of data at rest, in transit, and in use to protect against interception and unauthorized access.
- Leveraging AES-256 or higher standards, and managing keys securely through dedicated KMS platforms.
- Access control policies that enforce the principle of least privilege and prevent data misuse.

- Apply data anonymization, pseudonymization, and masking aligned to use-case risk and utility.
- Regular audits and monitoring of AI pipelines to detect anomalies or malicious activity.
- Dependency Risk Management: Review and monitor third party libraries and tools used in AI workflows to detect vulnerabilities or embedded threats.
- Adversarial training, where models are exposed to manipulated inputs during development to improve resilience. Measurable assurance techniques should be investigated (e.g., robust accuracy, provable defenses, and randomized smoothing).
- Secure deployment practices, such as authenticated API access, artifact signing for models, and watermarking for generated outputs and storing and fetching authentication credentials from a centralized, enterprise managed repositories as against using those in the codebase, to prevent unauthorized use.
- Non-Human Identities (NHI) like API key, token, certificates need to rotate periodically to maintain security.
- Zero Trust architecture and principals.

Equally important is the establishment of a governance framework that ensures transparency, compliance with regulations such as GDPR and California Consumer Privacy Act (CCPA), and ethical use of AI systems across their lifecycle.

# AI Controls Matrix Mapping

The CSA AI Controls Matrix (AICM) defines a comprehensive set of controls designed to manage data security and privacy across the AI lifecycle. Controls DSP-01 through DSP-24 focus on foundational principles, such as data classification, retention, protection, and privacy. These controls provide critical coverage for data-centric risk management in AI systems and align with applicable legal and regulatory requirements.

Each control supports one or more lifecycle stages including data collection, storage, processing, model training, deployment, and monitoring. Below is an overview of the AICM's current data security controls as they relate to AI environments.

## Data Security and Privacy Lifecycle Management Controls

### DSP-01: Security and Privacy Policy and Procedures
Establish and maintain data security and privacy policies that are reviewed at least annually. These policies must cover classification, handling, incident management and legal compliance across the data lifecycle.

### DSP-02: Secure Disposal

Apply secure disposal methods to prevent data recovery from decommissioned storage media.

### DSP-03: Data Inventory
Maintain a detailed inventory of sensitive and personal data assets.

### DSP-04: Data Classification
Classify data according to its type, sensitivity, and business impact to ensure appropriate protection and handling.

### DSP-05 Data Flow Documentation
Document where data is stored, processed, and transmitted. Review this documentation regularly, especially after significant system changes.

### DSP-06: Data Ownership and Stewardship
Assign and document data ownership and stewardship responsibilities, especially for personal and sensitive data.

### DSP-07: Data Protection by Design and Default
Ensure systems are built with default privacy protections that meet legal obligations.

### DSP-08: Data Privacy by Design and Default
Develop systems, products, and business practices based upon a principle of privacy by design and industry best practices. Ensure that systems' privacy settings are configured by default, according to all applicable laws and regulations.

### DSP-09: Data Protection Impact Assessment
Conduct impact assessments for systems processing personal data to evaluate privacy risks.

### DSP-10: Sensitive Data Transfer
Apply privacy and security measures to protect sensitive data in transit.

### DSP-11: Personal Data Access, Rectification and Erasure
Enable individuals to exercise rights over their personal data, including access, correction, and deletion.

### DSP-12: Limitation of Purpose in Personal Data Processing
Process personal data only for declared and lawful purposes.

### DSP-13: Personal Data Sub-processing
Apply safeguards and document how personal data is shared across the supply chain.

### DSP-14: Disclosure of Data Sub-processors
Disclose and evaluate sub-processor access to personal or sensitive data before allowing access.

### DSP-15: Limitation of Production Data Use
Restrict use of production data in non-production environments without appropriate authorization.

### DSP-16: Data Retention and Deletion
Ensure retention and deletion practices align with business needs and legal requirements.

### DSP-17: Sensitive Data Protection
Apply technical and organizational safeguards to protect sensitive data throughout its lifecycle.
- Challenge and document justifications for any use of sensitive data within the AI model.
- Segment sensitive data from non-sensitive transactional data.
- Monitor, control, and log all AI model access (including Agentic AI and NHI access) to all defined sensitive data.
- Sensitive data must be encrypted at rest, in transit, and in use.

### DSP-18: Disclosure Notification
Establish procedures to handle law enforcement requests for personal data in compliance with applicable laws.

### DSP-19: Data Location
Document where data is physically stored and processed, including backups.

### DSP-20: Data Provenance and Transparency
Track data sources and ensure transparency about data origin and use.

### DSP-21: Data Poisoning Prevention & Detection
Implement safeguards to detect and prevent data poisoning attacks in AI model training.

### DSP-22: Privacy Enhancing Technologies
Use PETs when processing training data, based on risk and business requirements.

### DSP-23: Data Integrity Check
Validate training and fine-tuning data for consistency and apply dataset versioning to ensure traceability.

### DSP-24: Data Differentiation and Relevance
Ensure training data is relevant to the AI model's intended use and properly differentiated.

### CEK-03: Data Encryption
Provide cryptographic protection to data at-rest and in-transit, using cryptographic libraries certified to approved standards.

The table below shows the different controls supporting various phases of the lifecycle.

| Data Collection | Data Storage | Data Preprocessing | Model Training | Model Deployment | Model Monitoring |
|---|---|---|---|---|---|
| DSP-03 | DSP-17 | DSP-05 | DSP-07 | DSP-15 | DSP-28 |
| DSP-04 | DSP-19 | DSP-22 | DSP-21 | DSP-12 | |
| | CEK-03 | | DSP-23 | | DSP-11 |

Refer to CSA AI Controls Matrix (AICM) Controls Mapping

# Control Gaps and Future Recommendations

While the current AICM provides strong coverage for core data protection, several emerging threats remain underrepresented or unaddressed. These include prompt injection, model inversion, federated learning governance, and unmonitored use of AI tools. To address these gaps, the following new controls are proposed for future inclusion:

- **DSP-25: Prompt Injection Defense:** Implements safeguards such as input validation, context isolation, and output filtering to prevent direct and indirect prompt injection attacks on LLMs.
- **DSP-26: Model Inversion & Membership Inference Protection:** Applies techniques, such as differential privacy, response truncation, and entropy monitoring, to limit model output leakage.
- **DSP-27: Federated Learning Governance:** Defines security and trust requirements for decentralized model training, including update validation and client authentication.
- **DSP-28: Shadow AI Detection**: Detects unauthorized use of AI tools and external LLMs using monitoring, DLP tools, and classification-based filtering.
- **Digital Identity Rights Framework (DIRF)**: In addition to DSP-25 through DSP-28, the Cloud Security Alliance has published the Digital Identity Rights Framework (DIRF) as a research blog artifact. DIRF introduces nine domains and 63 controls that directly address identity cloning, consent, memory drift, and monetization gaps currently unrepresented in the AI Controls Matrix.
- **Secure Runtime Environment**: Enforce a principle of least privilege for the AI application's runtime. This includes using hardened container images, isolating the environment from other network resources via segmentation, and implementing real-time monitoring of system calls and file access to detect anomalous behavior

These proposals reflect risks that are not explicitly addressed by existing AICM entries. A cross-domain analysis confirms that DSP-25 through DSP-28 are unique in scope and essential for covering GenAI-specific risks that extend beyond traditional data protection paradigms.

## Real-World Examples Supporting New Controls

- **DSP-25 (Prompt Injection Defense):** Aligned with [OWASP LLM Top 10 #1,](#) Unsanitized prompts have resulted in unauthorized system actions and information disclosure.
- **DSP-26 (Model Inversion & Membership Inference):** API leakage incidents, such as the [2023 GPT-4 exposure case,](#) highlight the risk of sensitive training data reconstruction.
- **DSP-27 (Federated Learning Governance):** Used in mobile and healthcare environments where decentralized model training requires robust coordination and trust models.
- **DSP-28 (Shadow AI Detection):** Prevents unmonitored use of GenAI tools that bypass internal data protection controls, such as confidential document uploads to public LLMs.

## Summary of AI Risk Control Gaps and Proposed AICM Extensions

| AI Risk | Mapped AICM Control(s) | Gap? | Proposed Control |
|---|---|---|---|
| Prompt Injection | None | Yes | DSP-25 |
| Model Inversion / Membership Inference | DSP-22 (partial) | Partial | DSP-26 |
| Federated Learning Governance | None | Yes | DSP-27 |
| Shadow AI / Unmonitored LLM Usage | None | Yes | DSP-28 |

These proposed extensions support CSA's GenAI Shared Responsibility Model and OWASP's LLM Top 10 guidance. They may serve as a foundation for future revisions of the AI Controls Matrix and inform peer review within CSA working groups and industry stakeholders.

# Data Security Risks Specific to AI

AI systems present a unique set of data security challenges due to the scale, sensitivity, and complexity of the data they process. These risks appear at every phase of the AI lifecycle and often fall outside the scope of traditional security frameworks. The following sections highlight core risk categories and emerging challenges tied to AI-driven environments.

# Inherent Risks

AI systems rely on vast volumes of heterogeneous data, much of it sensitive or proprietary. This dependency creates distinct vulnerabilities during data collection, storage, processing, and model execution.

## Risks in Data Collection and Storage

Centralized repositories, such as data lakes, are valuable targets for attackers and are prone to insider misuse without strong access controls.

- **Data Aggregation Risk:** Large-scale centralized storage increases exposure to breaches and data misuse.
- **Poor Source Vetting:** Unverified or open-source data can introduce poisoned content into the training pipeline.
- **Metadata Exposure:** Even seemingly harmless metadata, like logs or timestamps, can leak sensitive information if improperly secured.
- **Provenance/Origin loss**: Without clear records of where data originated, organizations cannot track data lineage or demonstrate regulatory compliance.

## Threats During Data Processing and Model Training

AI pipelines are complex and often span distributed infrastructure, increasing the attack surface at each stage.

- **Data Poisoning Attacks:** Malicious actors may inject compromised data that influences model behavior.
- **Training Infrastructure Compromise:** Unsecured environments or pipeline components can be exploited to extract raw data or interrupt training.
- **Insecure Intermediate States:** Temporary outputs such as embeddings or preprocessing artifacts may carry sensitive information and are often overlooked.

## Vulnerabilities in AI Algorithms and Output Data

Deployed AI models can leak sensitive data, act unpredictably, or reinforce bias when not properly controlled.

- **Model Inversion and Membership Inference:** Attackers may reconstruct or confirm the presence of training data using model outputs.
- **Overexposed Outputs:** Generative models can unintentionally reveal training data through poorly filtered responses or prompt manipulation.
- **Bias and Amplification:** Inadequately trained models may produce discriminatory outcomes, especially in sensitive areas like healthcare, finance, and hiring.

# Risk Categories

| Category | Risk Description | Examples |
|---|---|---|
| Data Privacy Risks | Unauthorized access or exposure of sensitive data | AI trained on PII without encryption or anonymization |
| Data Integrity Risks | Tampering or injection of adversarial data | Modified financial records influence model output |
| Model Inversion Attacks | Reconstructing sensitive data from model responses | Extracting medical history from a diagnostic model |
| Model Extraction & Theft | Cloning models via repeated API queries | Reverse-engineering a proprietary fraud detection model |
| Bias & Fairness Risks | Unintended discrimination or unfair outcomes | Hiring models favoring specific demographics |
| Data Leakage & Membership Inference | Identifying whether data was used in training | Confirming a name in a customer sentiment model |
| Adversarial Attacks | Manipulating inputs to produce false outputs | Altered images fool facial recognition systems |
| AI Supply Chain Vulnerabilities | Compromised or unverified third-party components in AI systems | Backdoors in pre-trained NLP models |

# Data and Model Risks of Being Outdated

AI models are only as effective as the data they rely on. Over time, both training data and model logic can become outdated, leading to reduced performance, misalignment with policy or user intent, and potential security gaps.

## Data Risks

- **Risks to Training Data:** Models trained on obsolete or biased data sets may produce inaccurate or harmful predictions.
- **Data Lineage:** Without source-to-use traceability, it becomes difficult to audit and verify data integrity.
- **Data Provenance:** Unknown or unverified sources increase the risk of misinformation, copyright violations, or embedded backdoors.
- **Data Quality:** Poor labeling, duplication, or noise can degrade model performance and introduce systemic vulnerabilities.

## Model Risks

- **Model Objective Drift:** Changing organizational or regulatory goals can render older models noncompliant or ineffective.
- **Model Algorithm Risk:** Legacy models may lack resilience against new adversarial techniques or include unpatched flaws.
- **Model Testing Gaps:** Outdated or infrequent testing reduces visibility into model accuracy and fairness.
- **Model Versioning and Reuse:** Deploying models without reassessment can reintroduce deprecated or vulnerable configurations.
- **Choice of Vulnerable Models**: Usage of models with open vulnerabilities or lack of the model vetting process can cause significant disruptions.

## Process Risks

AI development involves multiple systems, stakeholders, and workflows. Process-level misalignments can introduce hidden vulnerabilities.

- **Lack of Continuous Monitoring:** Without automation and real-time oversight, deviations in pipeline behavior may go unnoticed.

- **Over-reliance on Static Policies:** Outdated approval processes and access controls may expose models or data to unauthorized users.

# Delivery Risks

Deploying AI introduces another layer of complexity, especially across hybrid or edge environments.

- **Outdated Environments:** Models deployed on unsupported runtimes or insecure containers may introduce new threats.
- **Model Drift Post-Deployment:** Without retraining or feedback loops, model performance may degrade silently.
- **Dependency Insecurity:** Reliance on unmaintained third-party libraries or model hubs can introduce supply chain vulnerabilities.

# Recent Attacks on AI Systems

High-profile AI security incidents provide important lessons.

- **DeepSeek Database Exposure (2025):** Highlighted the risks of insufficient infrastructure monitoring and lack of access controls.
- **OpenAI API Incident (2023):** Demonstrated how improper session management can expose user-generated prompts and outputs.
- **DeepMind Retention Failures (2022):** Showed the difficulty of enforcing retention policies when data is reused for extended research purposes.

## Lessons learned from recent incidents

From these events, several best practices emerge:

- Enforce strong data lineage, model version control and granular access control.
- Continuously monitor deployed models for drift, bias, and anomalies.
- Implement secure session and API token management to protect user data in transit.
- Regularly review alignment with evolving regulatory frameworks and automate data retention.
- Apply change management protocols across the AI data pipeline.
- Integrate patching and update cycles into the broader cybersecurity strategy, embedding privacy-by-design principles throughout.

# Summary of Data Security Risks

AI systems introduce a layered and evolving risk landscape, including:

- Data-centric threats such as poisoning, unverified or low-quality sources, lineage ambiguity, and provenance gaps that increase exposure to misinformation, bias, and compliance violations.
- Model-level vulnerabilities including inversion, extraction, unfair bias, outdated algorithms, and insufficient adversarial robustness, which may compromise both performance and trust.
- Process and infrastructure risks such as insecure deployment environments, dependency on unmaintained libraries, weak monitoring, and static or outdated controls that fail to adapt to emerging threats.

To mitigate these risks, organizations must adopt a proactive security posture that emphasizes continuous assurance, lifecycle-aware governance, and security-by-design principles. Only by addressing each layer of the AI stack can organizations build resilient, secure, and trustworthy AI systems.

# Systematic Threat Taxonomy for AI Data Security

Understanding AI threats requires a structured approach based on lifecycle stage and security objectives. This taxonomy demonstrates the critical relationship between training-phase integrity failures and inference-phase confidentiality breaches.

## Training-Phase Integrity Attacks

Data Poisoning Classifications:

**Targeted Poisoning**: Attackers inject samples to cause specific misclassifications

- Example: Poisoning malware detector to classify ransomware as benign
- Detection: Statistical outlier analysis, data provenance verification
- Mitigation: Robust training with anomaly detection, validated data sources

**Backdoor Poisoning**: Hidden triggers embedded in training data create conditional vulnerabilities

- Example: Stop signs with specific patterns triggering misclassification
- Detection: Neural cleanse techniques, activation pattern analysis
- Mitigation: Input diversity validation, trigger pattern detection

**Label Manipulation**: Systematic corruption of training labels

- Example: Inverting security classifications in training datasets

- Detection: Cross-validation with independently labeled datasets
- Mitigation: Multi-annotator consensus, cryptographic label verification

## Inference-Phase Confidentiality Attacks

Privacy Extraction Methods:

**Model Inversion**: Reconstructing representative training samples from model outputs

- Primary Risk Factor: Models with excessive capacity relative to task complexity
- Defense: Differential privacy mechanisms with calibrated noise addition
- Monitoring: Analyze output confidence distributions for exploitation patterns

**Membership Inference**: Determining whether specific records existed in training data

- Primary Risk Factor: Model memorization due to overfitting
- Defense: Regularization techniques (dropout, weight decay, early stopping)
- Monitoring: Deploy shadow models to test for membership leakage

Vulnerability-Attack Correlation Framework

| Training Vulnerability | Associated Inference Attack | Relative Risk |
|---|---|---|
| Overfitting | Membership Inference | High |
| Poor Generalization | Model Inversion | Medium-High |
| Backdoor Presence | Targeted Data Extraction | Critical |
| Label Corruption | Unintended Information Disclosure | Medium |

# Prompt Guardrails

Establishing guardrails around AI prompts is essential to mitigate data leakage, enforce access control, and ensure regulatory compliance. As large language models (LLMs) become more deeply integrated into enterprise workflows, organizations must adopt a layered defense strategy that includes Data Loss Prevention (DLP), data classification, and AI-aware access controls.

# DLP

A robust Data Loss Prevention (DLP) solution is foundational for managing how sensitive data is accessed and used by AI platforms. If an organization already blocks the posting of confidential information to external websites, it may already have effective first-layer protections in place.

- **Context-aware DLP:** Modern DLP platforms can detect and block the transfer of regulated data types, such as personally identifiable information (PII), protected health information (PHI), or trade secrets to generative AI tools or unauthorized systems.
- **Prevention vs. Detection:** Proactive blocking of risky behavior is far more effective than retroactive investigation after a data loss event involving AI prompts.
- **SASE-powered network DLP**: As unmanaged applications and AI tools become more common, integrating data loss prevention with Secure Access Service Edge (SASE) enables inline discovery and control of sensitive data traversing the network addressing gaps beyond endpoint or managed app enforcement.
- **Zero-Trust Network Strategy**: A Zero Trust network strategy enhances DLP by continuously verifying users and devices, applying least-privileged access, and monitoring all data flows, even after initial authentication, to prevent sensitive data exfiltration from both external threats and insider risks

# Labeling Data

Effective data classification makes AI guardrails significantly easier to implement and enforce. Labeling data with sensitivity levels or project tags enables organizations to define automated rules for handling prompts.

- **Project-Based Controls:** If a confidential project (e.g., "Falcon") is appropriately tagged, systems can automatically prevent related data from being uploaded or referenced in AI interactions.
- **Regulated Data Types**: In sectors like healthcare, where PHI or PII is tightly regulated, enforcing labeling policies can significantly reduce risk exposure when AI tools are in use.
- **Integration with DLP:** Most enterprise DLP solutions support the ability to detect and act upon sensitivity labels, making this a scalable approach to policy enforcement.

# Advanced Techniques to Modulate Responses

**1. Prompt Engineering & Structured Prompts**:

- Carefully designed prompts can prevent the accidental disclosure of sensitive inputs (such as PII or confidential business data) to AI models.

**2. Reinforcement Learning with Human Feedback (RLHF)**:
- Include security-specific human feedback during training (flagging outputs that could reveal secrets).
- Multi-objective RLHF balancing factuality, safety, and confidentiality.

**3. Post-Processing / Response Filtering**:
- Use classifiers to detect PII, financial, or proprietary data in generated responses.
- Apply context-aware sanitisation or anonymisation before release.

**4. Controlled Generation Techniques**:
- Biasing token probabilities to avoid sensitive terminology or restricted topics.
- Use top-k/top-p sampling to control randomness, reducing accidental exposure.

**5. Model Modularisation / Ensemble Techniques**:
- Assign a security-focused module to vet outputs before release.
- Isolate models handling confidential datasets from general-purpose models.

**6. Embedding & Retrieval-Augmented Generation (RAG)**:
- Retrieve from encrypted, access-controlled repositories.
- Embed policies or sanitised data as context to modulate outputs.

# Guardrails Enable AI Adoption

AI guardrails are not just defensive mechanisms, they are strategic tools that support responsible AI adoption.

- **Safe Enablement vs. Blanket Restriction:** Attempting to ban AI outright through policy is unlikely to be effective and may drive employees to use unvetted tools via shadow IT.
- **Cultural Acceptance and Trust:** By transparently deploying controls that protect sensitive data while allowing approved use, organizations can foster a culture of trust and accountability. Provide user training to raise awareness of how users can support the safe and secure use of AI tools.
- **Acceleration of Innovation:** Guardrails help teams safely explore AI capabilities, accelerating adoption and maximizing business value.
- **Continuous Learning and Feedback**: Facilitate organized feedback mechanisms in AI applications, allowing teams to identify errors, mitigate biases, and enhance model performance progressively. This approach guarantees that AI development proceeds securely and efficiently, consistently aligning with organizational objectives and ethical principles.

- **Regulatory and Compliance Readiness**: Include mechanisms and oversight that adhere to regulatory requirements, helping organizations to demonstrate accountability and compliance proactively. This reduces risks during audits and enhances confidence in AI systems, facilitating faster adoption.

## Online Tokenization

Online tokenization is a policy-enforced runtime control that prevents sensitive data from entering model contexts and enables restoration only under explicit authorization. It operationalizes privacy by protecting prompts and retrieved context while preserving utility and maintaining clear authorization and audit boundaries.

- **Scope:** Real-time prompts, retrieved context, and model outputs across chat, tools, and API calls; streaming-safe over SSE/gRPC without buffering entire payloads.
- **Enforcement:** Pre-request scan + tokenization; response-only detokenization gated by RBAC/ABAC with reason codes; token domains scoped to session/tenant with no cross-client reuse.
- **Telemetry:** Immutable audit of who/what/when/why for every detokenization; detector confidence tracking; gateway latency with attention to the slowest 1–5% of requests; counts of tokenized items; approved vs. denied detokenization events.
- **Impact:** Mitigates accidental leakage, impedes prompt injection-driven exfiltration of secrets, and enables safer GenAI adoption without relying solely on post-response filters.
- **Fallbacks:** On low detector confidence, block or escalate (fail-closed); apply circuit breakers and graceful degradation paths with clear caller messaging.

# AI Ethical and Legal Considerations

Artificial intelligence introduces complex legal and ethical challenges that impact everything from intellectual property to human rights. These concerns are not theoretical; they are already shaping regulation, public trust, and litigation. As AI systems become more autonomous and influential in decision-making, organizations must embed accountability, consent, and transparency into their AI development and deployment practices.

## Legal Issues

**Intellectual Property:**

Content created by AI raises unresolved questions about intellectual property law. In some jurisdictions, ownership may vary depending on whether the data, model, or output is being considered. Rights may apply through patents, trade secrets, copyright, or database protections. Each situation requires legal review based on how the model was trained and how its outputs are used.

**Copyrighted material:**
Meta has been involved in ongoing legal action after reportedly using 82 terabytes of pirated media to train an AI model. This case illustrates the potential consequences of unauthorized data use in AI development.

**Libel and Defamation:**
When an AI system generates content that turns out to be defamatory, accountability remains unclear. Legal responsibility may fall on the developer, deployer, or user, depending on the context and applicable laws.

**Data poisoning:**
Large language models are often trained on publicly available data. This opens the door for malicious actors to inject false information at scale in an attempt to influence future AI outputs.

**Cross-border Compliance and Jurisdictional Ambiguity:**
AI systems deployed globally must navigate overlapping and sometimes conflicting regulations such as the GDPR and CCPA. Determining which laws apply becomes difficult when models are trained in one region and used in another.

**Accountability and Liability:**
Legal responsibility for harm caused by AI, such as financial fraud or medical errors, is still evolving. While AI liability is evolving, frameworks such as the EU AI Act are beginning to define responsibility. These regulations often differentiate between the roles of the AI model creator and the deployer, assigning liability based on the risk level of the system and the actions of the parties involved.

**Contractual Risk in AI Vendor Ecosystems:**

- Vendor contracts often lack provisions for AI-specific concerns, such as data governance, auditability, and breach response. This creates gaps in shared liability and complicates incident resolution.
- CSA's Digital Identity Rights Framework (DIRF), published as a research blog artifact, further strengthens these areas by requiring explicit user consent, transparency in identity usage, and auditability of clone governance. DIRF introduces enforceable safeguards that extend beyond current legal interpretations, helping organizations operationalize consent and digital identity

rights in AI environments.

# Ethical Considerations

**Steganography in AI-Generated Media:**
AI-generated images may contain hidden data, either by design or as a byproduct of model behavior. These embedded elements can be used maliciously, raising questions about the trustworthiness of synthetic media.

**Deceitful Behavior in AI Systems:**
Some AI systems have demonstrated the ability to deceive under certain conditions. Whether intentional or emergent, this behavior undermines trust and supports the need for strict limitations on data access.

**Hallucinations (factually incorrect or nonsensical output by an AI system) and Unreliable Outputs:**
AI models can produce false or misleading results, especially when dealing with unfamiliar or ambiguous prompts. For example, many AI image generators consistently depict watch faces set to 10:10, reflecting the bias of their training data. Efforts to retrain models to correct this behavior have had limited success, raising concerns about the reliability of AI in even simple tasks.

**Consent and Ownership in Training Data:**
Models are frequently trained on scraped public data without consent from the original authors. This includes personal blogs, health forums, and other sensitive sources, leading to ethical questions around digital consent and data ownership.

**Bias Amplification and Disparate Impact:**
Even when not explicitly biased, AI systems may reinforce existing disparities present in their training data. This is particularly harmful in areas such as hiring, lending, and healthcare, where affected individuals may already face systemic disadvantages.

**Synthetic Media and Deepfakes:**
AI-generated content can be nearly indistinguishable from real media. Without trusted ways to verify authenticity, such as provenance metadata or embedded watermarks, these tools can be used to spread misinformation or commit fraud.

**Emergent Behavior and Misuse Potential:**
Advanced AI systems may display unexpected capabilities such as manipulation, deception, or autonomous planning. Ethical responsibility requires developers to anticipate and mitigate these risks during design and testing, not just after deployment.

The ethical and legal risks posed by AI are not hypothetical. They are already shaping regulatory frameworks, court cases, and public expectations. Organizations must address these risks proactively by integrating legal reviews, transparent data practices, and ethical oversight into every phase of AI development. Doing so builds public trust and supports the safe and responsible use of AI technologies.

# Regulatory and Compliance Landscape

AI systems are increasingly subject to a wide range of regulations and compliance frameworks that govern how personal and sensitive data is collected, processed, stored, and shared. This is particularly critical in environments that involve automated decision-making or high-risk data usage. The following section provides an overview of key regulations, frameworks, and emerging standards shaping the secure and responsible development of AI.

## Key Regulations

### United States and Global Overview

More than 450 AI-related bills have been introduced across U.S. states. Similar efforts are underway in Canada, Brazil, Japan, and other countries. The European Union's AI Act (In force as of 2024 First act effective Feb2025, full compliance 2027) represents the most comprehensive attempt to establish risk-based compliance obligations for AI systems, including mandatory requirements for transparency, bias testing, and data governance in high-risk use cases.

### General Data Protection Regulation (GDPR)

A foundational privacy regulation for AI systems operating in or processing data from the EU.

## California Consumer Privacy Act (CCPA)

Applies to entities collecting personal data from California residents.

- **AI-Specific Provisions:** Requires explicit notice of AI usage, transparency in profiling, and the right to opt out of automated decision-making.
- **Technical Requirements:** Organizations must maintain a detailed inventory of training datasets and document how personal data is used within AI systems.
- **Consumer Rights:** Grants individuals access to, and deletion of, personal data used for AI model training and inference.

## Health Insurance Portability and Accountability Act (HIPAA)

Regulates the use of Protected Health Information (PHI) in AI-driven healthcare systems.

- **Security Safeguards**: Enforces encryption, access controls, and audit logging for PHI used in AI training or inference.
- **Vendor Compliance**: Business Associate Agreements (BAAs) must include AI vendors, and AI-specific risks must be addressed in security assessments.

## Asia-Pacific AI Regulations

- **China's AI Regulations:** provisions for algorithmic recommenders, regulation of deep synthesis
- **Singapore's Model AI Governance Framework:** a self-assessment toolkit for deploying responsible AI
- **Japan's AI Principles:** Society 5.0 and human-centric AI
- **India's Responsible AI Framework:** fairness, accountability, and transparency

## Key AI-Specific Requirements

### Algorithmic Transparency

Requires meaningful explanation of automated decisions and how model logic affects outcomes.

**Data Minimization**

Limits data collection to only what is necessary for the specific AI function and mandates regular review and deletion protocols.

**Data Subject Rights**

Includes the right to opt out of automated decisions, request human intervention, and contest AI-driven outcomes.

# Technical Frameworks and Standards

## NIST AI Risk Management Framework (RMF)

Provides voluntary, widely adopted guidance on how to identify, assess, and mitigate risks associated with AI development. Focus areas include trustworthiness, socio-technical risk, data integrity, and explainability.

## ISO 23894 and ISO/IEC 42001

- **ISO 23894:** Risk management for AI systems, focused on integrating risk controls into design and deployment workflows.
- **ISO/IEC 42001:** First formal standard for AI governance management systems, enabling organizations to implement policies for ethical and secure AI operations.

## Cloud Security Alliance (CSA) Guidance

- **AI Controls Matrix (AICM):** A detailed set of technical and organizational controls mapped to the AI lifecycle. Future extensions will include controls for threats like prompt injection, model misuse, and federated learning governance.
- **Security, Trust, Assurance and Risk (STAR) Program (Level 3):** Expands on existing STAR Levels 1 (self-assessment) and 2 (audited) by introducing AI-specific control validation, maturity scoring, and evidence-based attestation using automation.
- **Compliance Automation Revolution (CAR):** A next-generation CSA initiative that builds upon STAR Level 3 to automate AI-related compliance workflows. CAR uses AI and tooling to evaluate risk attestations for consistency, traceability, and maturity. This approach accelerates trust and transparency in cloud-based AI deployments.
- **Generative AI Shared Responsibility Model:** Provides a clear delineation of roles and responsibilities between cloud providers, AI vendors, and customers. Helps organizations identify

where to apply specific controls across the AI stack.

## OWASP LLM Top 10

Provides voluntary, widely adopted guidance on how to identify, assess, and mitigate risks associated with AI development.

## MITRE ATLAS Framework

Catalogs real-world adversarial techniques targeting AI systems. It supports threat modeling, red teaming, and the development of defensive controls.

## ENISA Guidelines for Secure AI Development

Issued by the EU Agency for Cybersecurity, these guidelines include secure development principles, adversarial robustness, and data pipeline protections.

## IAPP AI Governance Framework

Focuses on aligning data privacy practices and governance protocols with AI operational needs and ethical requirements

## OECD AI Principles

Adopted by over 45 countries, these high-level guidelines promote human-centric, transparent, and accountable AI development.

## UNESCO AI Ethics Recommendations

The first global normative framework focused on ethical AI. It addresses issues such as algorithmic bias, cultural sensitivity, and data stewardship.

## EU AI Act (Pending Adoption)

- Classifies AI systems into four tiers: unacceptable, high risk, limited risk, and minimal risk.
- Imposes conformity assessments, data quality requirements, and mandatory human oversight for high-risk systems.

# Sector-Specific Frameworks

### IMDRF AI Guidance (Healthcare)

Guidelines for medical AI and machine learning models from the International Medical Device Regulators Forum. Covers real-world model monitoring, transparency, and version control.

### FINRA AI Guidelines (Finance)

Regulatory expectations for U.S. financial institutions using AI. Focus areas include explainability, fairness, and robust audit trails.

### NAIC Model Law on Insurance AI

A model law for insurers, developed by the National Association of Insurance Commissioners. Requires documentation of AI usage, fairness evaluation, and consumer rights protections.

### FATF Guidance on AI and Digital Identity

Published by the Financial Action Task Force, this guidance focuses on using AI responsibly in anti-money laundering systems and digital identity verification.

| Framework | Focus Area | Relevance to AI Data Security |
|---|---|---|
| GDPR | Privacy Regulation | Legal basis for using personal data in AI systems |
| ISO/IEC 42001 | AI Governance | AI management system standard |
| NIST AI RMF | Risk Management | Framework for identifying and mitigating AI risks |
| CSA AI Controls Matrix (AICM) | Cloud AI Controls | Mapped technical and organizational safeguards for AI |
| ENISA Secure AI Guidelines | Cybersecurity | Best practices for secure AI development and deployment |

| OECD AI Principles | Policy and Ethics | Global guidance on transparency, fairness, and accountability |
|---|---|---|

## AI Shared Responsibility Model

The shared responsibility model for AI clarifies how roles and security obligations are distributed across cloud providers, AI platform vendors, and enterprise users. It highlights where each party must apply controls to secure data and maintain compliance across AI development and deployment environments. Align to the AI Organizational Roles and Responsibilities and maintain consistency in reference and avoid standard exhaustion by referring to the standards in play vs. repeating content.



Comparative shared responsibility

# Privacy Enhancing Technology for AI

## Data Minimization

Only collect and process data that is strictly necessary for the AI task.

# Data Classification

Classifying data by type and sensitivity is foundational for governing and securing AI applications. Data types generally include:

- **Structured:** Databases, spreadsheets
- **Unstructured:** Emails, PDFs, text documents
- **Semi-structured:** JSON, XML files

Each data asset should be labeled according to sensitivity, with aligned protections:

- **Public**
- **Company Internal**
- **Confidential**
- **Restricted**

# Data Anonymization

Anonymization removes identifiers entirely, while pseudonymization replaces them with fictional values. These processes enable privacy-preserving analysis while obscuring personal identifiers.

Effective anonymization requires thorough data discovery and classification to locate sensitive fields. Common techniques include:

- **Redaction:** Removing fields completely
- **Generalization:** Replacing with broader categories
- **K-anonymity:** Ensuring records blend into groups
- **Noise Injection:** Altering data slightly to mask identity

Built-in tools are available from:

- **Database providers:** Oracle, SQL Server, MySQL
- **Cloud platforms:** AWS, Azure, GCP

# Data Masking

Data masking creates fake versions of sensitive datasets that retain their structure but hide actual values. This enables safe AI training and development without exposing real data. Common masking tools include:

- IBM InfoSphere Data Masking
- Informatica PowerCenter
- Precise DataFlex
- Broadcom Test Data Manager
- Denodo Data Masking (open-source)

# Encryption

**Homomorphic Encryption:** Allows computations directly on encrypted data, preserving privacy even during AI model use. While computationally intensive, it ensures strong privacy guarantees.

Tools and platforms include:

- **OpenFHE:** Open-source HE library
- **Microsoft SEAL:** SEAL library for encrypted computation
- **IBM Cloud HEaaS (Homomorphic Encryption as a Service):** Cloud-based HE development
- **Enveil:** Commercial HE solutions for AI and analytics

**Searchable Symmetric Encryption (SSE):** SSE enables keyword searches over encrypted data sets without needing full decryption. This is useful for applications, such as medical records, financial transactions, or legal archives.

SSE capabilities are offered by:

- **Google Cloud and Amazon S3:** Through searchable index features
- **Paperclip SAFE**: Crypto-agile, SaaS encryption-in-use platform
- **PySearchable**: Open-source SSE Python library

# Tokenization of Data

Tokenization replaces sensitive values (PII/PHI/secrets) with secure representations (tokens) so that models, pipelines, and logs can operate without seeing the original data.

Unlike encryption (reversible with a key) and hashing (irreversible), tokenization can be reversible under strict control (via a vault/service) and deterministic when you need stable joins and auditability. It complements anonymization and masking and directly advances **DSP-17 (Sensitive Data Protection)** and **DSP-22 (Privacy Enhancing Technologies)**, while reinforcing **DSP-10 (Sensitive Data Transfer)**, **DSP-16 (Data Retention & Deletion)**, and **DSP-20 (Data Provenance & Transparency)** in the AICM.

## Design patterns you can reuse

- **Vaulted vs. vaultless**
  - *Vaulted*: a secured mapping table "value → token" lives in a tokenization service with KMS/HSM-backed keys, fine-grained access, and audit trails.
  - *Vaultless*: the token is derived deterministically (e.g., via FPE/HMAC) without a central lookup table.
- **Deterministic vs. random**
  Deterministic tokens enable deduplication, joins, and consistent analytics. Random tokens reduce enumerability and correlation risk.
- **Format-preserving vs. placeholders**
  *Format-preserving* (e.g., for numbers, IDs) keeps schemas happy. *Placeholders* (e.g., `[[EMAIL]]`, `[[ID]]`) are clearer to humans/LLMs and reduce accidental leakage.
- **Token domains/namespaces**
  Separate domains per data type (email, document ID, account) and per tenant to avoid cross-correlation.
- **Lifecycle management**
  Treat tokens and mappings as sensitive artifacts with retention windows, rotation, and break-glass procedures.

## Batch tokenization

**Applies to:** data at rest, training preparation, and retrieval-augmented generation (RAG).

- **Discovery & classification:** Inventory PII/PHI/secrets across lakes and warehouses; feed findings into policies.
- **Policy per data class:** Choose token types: FPE for rigid numeric fields; placeholders for free text; deterministic tokens where joins matter.
- **Determinism & idempotence:** The same input should always yield the same token to keep pipelines consistent.
- **Key & vault governance:** Isolate the tokenization service; manage keys in KMS/HSM; log every (de)tokenization event.
- **Quality controls:** Sample data to check detection recall/precision; monitor collisions; validate format constraints.
- **Pipeline placement:** Tokenize early (bronze → silver → gold) so derived sets are clean by construction; avoid leaking raw values into temp tables and job logs.
- **Embeddings & RAG:** Prefer tokenizing before creating embeddings so vector stores don't memorize PII. Where needed, use standardized placeholders that won't distort similarity search.
- **Data subject rights:** Keep a secure, auditable link from token to subject to support access/rectification/deletion without re-exposing raw data.

## Online tokenization

**Applies to:** real-time prompts, retrieved context, and model outputs.

- **Prompt gateway:** Place a proxy in front of LLMs. It detects and tokenizes sensitive data before sending to the model and detokenizes on the way back only if the caller is authorized.
- **Session-scoped tokens:** Issue ephemeral tokens per session/tenant; don't reuse mappings across clients.
- **Authorization to detokenize:** Enforce ABAC/RBAC plus "reason codes"; require justification and log who detokenized what, when, and why.
- **Streaming-safe:** Support SSE/gRPC streaming without buffering entire payloads; tokenize in-flight.
- **RAG connectors:** If detokenization is absolutely necessary, do it in the trusted source connector, not inside the LLM context.
- **Fail-safe defaults:** When the PII detector is uncertain, prefer *block or escalate* to avoid silent leakage.

## Patterns of Tokenization Implementation

- **Vaulted tokenization**: Centralized token vault backed by HSM–backed keys moves beyond FIPS 140-2 Level 3 (High security level).
- **Format Preserving Tokenization (FPT)**: It maintains the format of the tokenized value putting in place the sensitive value. This has real advantages for legacy applications, some of which may not translate sensitive values well and/or do not validate sensitive value forms.
- **Contextual tokenization**: Tokens are distinct for the same value based on the context that it is being used to mitigate correlation attacks.
- **Ephemeral tokenization**: Tokens are tied to session and expire after time windows or transactions.

## Trade-offs and limitations

- **Inference risk remains:** Tokenization hides direct identifiers but doesn't stop attribute inference from context; pair with anonymization or differential privacy where needed.
- **Operational surface:** Tokens can leak through logs, retries, and dead-letter queues, if your proxy scope is incomplete.
- **Latency:** An additional tokenization hop elevates the 95th/99th-percentile response times; employ scale-out, circuit breakers, and explicit fallback modes to contain the impact.
- **Semantic impact:** Placeholders help privacy but can reduce LLM utility; plan selective, **authorized** detokenization in post-processing, not in-prompt.

## What success looks like

- **Clean traffic:** Share of prompts/contexts reaching the LLM with no raw PII/PHI/secrets after the gateway. *SLI:* sanitized ÷ total. *Example target:* ≥ 99.5% (monthly).
- **Tail latency:** End-to-end latency across tokenization/detokenization, focusing on the slowest 5% and 1% of requests. SLI: share of requests completing within thresholds (ingress→egress, in-region). Example target: 95% ≤ 120 ms; 99% ≤ 250 ms.
- **Detector quality:** False-negative rate from periodic, stratified sampling/annotation (include red-team prompts). *Example target:* ≤ 1% per month.
- **Detokenization governance:** 100% of detokenizations authorized and justified; track approved vs. denied by role/team.
- **Leakage trend:** Severity-weighted leakage incidents before/after rollout (quarterly), aiming for a sustained downward trend.

# Data Resiliency

**Secure Multi-Party Computation (SMPC):** Allows multiple institutions to analyze data collaboratively without revealing their individual datasets. This enhances privacy across distributed AI projects.

Notable SMPC platforms and tools include:

- Cloud-ready platforms: Azure, AWS, GCP
- Ironclad Systems
- Sepior
- Duality Technologies
- OpenSMPC
- MP-SPDZ
- TF Enclave (TensorFlow)

**Federated Learning:** Enables AI training on user devices without transmitting raw data. Only model updates are sent back, preserving local privacy.

Federated learning is widely applied in:

- Mobile applications
- Healthcare AI
- IoT and smart devices

# Additional Privacy-Preserving Techniques

- **Local Differential Privacy:** Adds noise directly on user devices before data is used in model training, enabling private federated learning.
- **Privacy-Preserving SVM and KNN Models (Support Vector Machine (SVM) and K Nearest Neighbours (KNN)):** Enables classification and regression tasks without exposing training data.
- **Differential Privacy:** A mathematical framework ensuring that insights learned from datasets cannot be traced back to individuals.
- **Supporting tools include:**
  - TensorFlow Privacy
  - OpenDP
  - Research implementations from Apple, Google, and Microsoft
- **Transparent and Explainable AI (XAI):** XAI improves trust by enabling users to understand AI model decision-making. This is vital for identifying privacy risks and ensuring regulatory compliance.
  - Popular XAI toolkits include:
    - **TF-XAI** – TensorFlow's explainability framework
    - **DARPA XAI Toolkit** – Tools developed under DARPA's Explainable AI program

# Secure Data Storage and Transmission

Securing how data is stored and transmitted is critical to protecting sensitive information and ensuring the trustworthiness of AI systems. These protections must apply throughout the entire AI lifecycle, including data ingestion, training, testing, deployment, and inference. A strong security foundation begins with access controls, encryption, and ensuring data is properly segmented and classified. When these measures are lacking, the result can be data breaches, stolen models, regulatory violations, and degraded AI performance.

# Access Control Mechanisms

Strict access management is essential to prevent unauthorized interaction with sensitive data and models.

- **Role-based access control (RBAC)** limits access to only those users whose roles require it, helping enforce the principle of least privilege.
- **Attribute-based access control (ABAC)** evaluates factors like device security, user location, and time of access to provide context-aware policy enforcement.
- **Just-in-time (JIT) access** temporarily grants permissions only when needed, reducing long-term exposure to sensitive environments.

- **Audit logs** track who accessed what data and when, supporting accountability, compliance, and post-incident investigations.

# Encryption at Rest and in Transit

Encryption is essential to safeguard data from interception or tampering during storage and transmission.

- **At rest**, encryption protocols such as AES-256 should be applied to all files and data sets, including training data, model weights, and logs.
- **In transit**, communications between systems must use secure transport layers like TLS 1.3 to prevent eavesdropping or man-in-the-middle attacks.
- **Key management systems** like AWS KMS, Azure Key Vault, GCP Secrets Manager or HashiCorp Vault should be used to securely generate, rotate, and store encryption keys, with strict access controls and audit trails.

# Data Partitioning and Isolation

- **Environment isolation** ensures that development, testing, staging, and production systems do not share data or models unless explicitly required.
- **Tenant isolation** is important in multi-tenant AI platforms and should be enforced using containerization, virtualization, or network segmentation.
- **Tokenization** replaces sensitive values with secure placeholders, allowing safe use of representative data without exposing real information.
- **Trusted execution environments (TEEs)** such as Intel SGX or ARM TrustZone can protect in-use data and model inference by running processes in isolated, hardware-enforced memory regions.
- **Tokenization** should be used as a containment layer alongside environment and tenant isolation. Sensitive values are replaced with tokens before leaving the source domain; applications and models consume only tokens, while detokenization occurs conditionally inside a trusted service (vault) with KMS (or HSM) managed keys, per-request authorization, and immutable audit trails (logs). This reduces blast radius, enforces least privilege, and decouples real data from AI runtime layers. Where in-use access to original values is unavoidable, combined with Trusted Execution Environments (e.g., Intel SGX, ARM TrustZone) so any detokenization occurs inside hardware-isolated memory and does not escape the enclave.

# Secure Model Hosting and API Protection

Deployed models must be protected against unauthorized use, probing, and extraction.

- **API authentication and rate limiting** help control who can access model endpoints and how often, reducing the risk of scraping or brute-force attacks.
- **Output filtering** can detect and suppress sensitive content that may be inadvertently produced by models, especially in generative AI systems.
- **Watermarking and fingerprinting techniques** embed hidden identifiers in model artifacts or outputs to help detect and trace unauthorized distribution or misuse.

# Leveraging AI for Data Security and its Implications

## Using AI to Enhance Data Security Systems

The need to secure and continuously safeguard critical information technology assets has always been a major concern for countries and businesses. This imperative has become even more severe with the growing level of cyber attack (and their sophistication) faced by nearly all human organizations that use IT systems with varying levels of resilience across organizations and national economies. The World Economic Outlook's 2024 Global Cybersecurity Outlook (GCO) Report shows that there is a widening gap in distance between organizations sufficiently resilient to cyberattacks and those that are not thus leaving a wide and empty middle-resilient organizations[1]. The cost and speed of adoption of modern and more up-to-date technologies for ensuring adequate cybersecurity have been identified as the key factors driving this divide. While there is no guarantee that this gap will be closed anytime soon, the application of AI might provide a leeway to helping small organizations leapfrog their cyber resilience.

AI can be applied to enhance data security by leveraging its ability to learn, reason, solve problems and adapt to new and continuously changing situations. A key advantage that AI has over traditional security systems lies in its ability to automate complex processes, analyze a vast amount of data, and detect and respond to evolving and complex cyber threats. AI uses techniques including machine learning, natural language processing, and pattern recognition technologies to both secure and improve the resilience of IT systems to cyber attacks. The common objective of AI systems is to replicate human-like cognitive capabilities at scale to help proactively identify and resolve risks and potential threats to IT systems.

---

[1] 2024 Global Cybersecurity Outlook (GCO) available at
https://www3.weforum.org/docs/WEF_Global_Cybersecurity_Outlook_2024.pdf

# AI-driven Threat Detection and Response

AI enables more intelligent and adaptive approaches to identifying and responding to threats. Key techniques include anomaly detection, behavior profiling, and predictive analytics.

## Machine Learning for Threat Detection

Machine learning algorithms are a cornerstone of modern threat detection. These systems are trained on historical data to understand what constitutes normal behavior for users, devices, and applications. Once trained, they monitor systems in real time and flag any activity that deviates from these patterns.

- **Anomaly Detection**: AI systems use historical baselines to detect deviations from expected behavior. These deviations may indicate attempted intrusions, malware activity, or insider threats. The system continuously updates its understanding as behaviors change over time.
- **Behavioral Modeling**: Machine learning models analyze user activity, application usage, and device interactions across an organization. This allows them to identify subtle or novel threats that rule-based systems often miss.
- **Contextual Understanding**: AI systems can incorporate natural language processing to analyze unstructured sources such as log files, alerts, and threat intelligence feeds. This enhances situational awareness and supports smarter decisions.

## Predictive Analytics for AI-driven Threat Detection

AI also powers predictive analytics, which helps organizations stay ahead of threats rather than simply reacting to them.

- **Big Data Analytics**: AI-driven platforms process enormous volumes of real-time security telemetry. These tools aggregate information from endpoints, networks, and cloud environments to identify suspicious behavior and prioritize alerts.
- **Threat Hunting**: Using data-supported models, analysts can proactively search for indicators of compromise based on patterns identified in similar incidents across other environments.
- **Risk Prediction**: Statistical modeling and pattern recognition are used to anticipate potential threats. By analyzing system behavior, user activity, and threat intelligence, predictive analytics tools identify early warning signs of compromise.

These predictive capabilities help reduce alert fatigue, improve detection speed, and support more effective responses. They also play a crucial role in prioritizing risks and allocating security resources efficiently.

## AI-Assisted Data Security Operations (DSO)

AI-assisted DSO converts detections into operational controls by automating classification, tokenization, redaction, and policy enforcement, enabling organizations to prevent data exposure rather than respond after incidents.

- **Applies to:** Security operations runbooks and AI data pipelines, covering ingestion, storage, training, inference, and monitoring.
- **Detection to enforcement:** Orchestrate controls from detections: automatically classify sensitive fields; tokenize or redact before storage and use; enforce least-privilege at access; and trigger retention and deletion workflows. *(AICM: DSP-17, DSP-22, DSP-16)*
- **Privacy-preserving telemetry:** Treat logs, traces, prompts, and retrieved context as sensitive. Apply context-aware classification at ingestion; perform online tokenization or field-level redaction; and keep conditional detokenization within a trusted vault backed by KMS or HSM, with lineage and provenance recorded for audit. *(AICM: DSP-10, DSP-20, DSP-23)*
- **Action guardrails and human-in-the-loop:** Require RBAC and ABAC checks, dual control for high-risk actions, reason codes, and immutable audit trails. LLM and SOC copilots may suggest remediation but do not execute without policy-backed approval. *(Links to Prompt Guardrails; Sensitive Data Protection)*
- **Assurance and model health:** Set measurable quality and performance targets (precision and recall; end-to-end latency across tokenization and gateway paths, focusing on the slowest one to five percent of requests). Monitor drift and false-negative rates; conduct red-team exercises aligned with ATLAS and OWASP; and gate releases through dataset and model versioning and change control. *(AICM: DSP-23, DSP-20)*

# Benefits and Limitations of AI in Cybersecurity

## Benefits

- **Speed and Scale**: AI can process data far faster than human analysts and scale across large, complex environments.
  **Real-Time Detection**: Machine learning systems can respond to anomalies as they occur, improving incident response times.
- **Continuous Learning**: AI models improve over time by learning from new data and adapting to changing attack patterns.
- **Cost Efficiency**: Automation of repetitive tasks frees up human analysts for higher-value work and may lower overall operational costs.

## Limitations

- **Bias and False Positives**: AI systems may inherit biases from training data or produce inaccurate alerts if not properly calibrated.
- **Data Dependency**: The accuracy of AI depends on the quality, diversity, and volume of data available for training.
- **Complexity**: Implementing and managing AI solutions requires specialized expertise and integration with existing tools.
- **Adversarial Attacks**: AI models themselves can be targeted through techniques, such as data poisoning or adversarial inputs.

# Case Studies and Practical Applications

Real-world incidents provide critical insight into how data security risks emerge across the AI lifecycle. The following high-profile examples reveal weaknesses in access control, governance, and lifecycle protections. Each case reinforces the importance of applying AI-specific security measures.

## DeepSeek Database Exposure (2025)

- **Risk Type**: Access control failure and metadata exposure
- A publicly accessible AI infrastructure exposed a DeepSeek database, allowing access to internal log streams and sensitive configurations.
- **Lessons Learned**: AI systems must apply strong access management and infrastructure-level protections. Security controls should extend to metadata, logs, and ancillary data stores.

## Snowflake Cloud Data Breach (2024)

- **Risk Type**: Supply chain compromise and token reuse
- Attackers accessed sensitive customer data hosted by Snowflake by exploiting stolen credentials and weak token hygiene, affecting multiple enterprise clients including Ticketmaster.
- **Lessons Learned**: Third-party platform risks must be continuously monitored. API tokens, session keys, and credentials require lifecycle management and revocation policies. Zero trust access and session validation are critical to limiting blast radius from credential theft.

# OpenAI GPT-4 API Session Leakage (2023)

- **Risk Type**: Session isolation failure and prompt data leakage
- A flaw in API session handling allowed some users to see other users' chat histories.
- **Lessons Learned**: Secure API usage requires proper authentication, encryption, and session isolation. Prompt inputs and outputs should be treated as sensitive data.

# Google DeepMind Data Retention Non-Compliance (2022)

- **Risk Type**: Data lifecycle mismanagement
- Behavioral and medical datasets used in training were stored longer than allowed by applicable policies and regulations.
- **Lessons Learned**: Data retention and deletion must be governed throughout the AI pipeline. Models and datasets should be routinely audited for compliance.

# Apache Log4Shell Vulnerability (2021)

- **Risk Type:** Supply chain exposure and insecure dependency
- A critical zero-day vulnerability in the widely used Log4j Java logging library allowed attackers to execute arbitrary code on affected systems. Many AI infrastructure components and data processing tools that relied on Log4j were vulnerable to this exploit.
- **Lessons Learned:** Dependencies in AI pipelines must be continuously scanned and updated. Supply chain risks require SBOM (Software Bill of Materials), automated patching workflows, and real-time exposure monitoring.

# Clearview AI Facial Recognition Breach (2020)

- **Risk Type**: Unsecured training data and third-party access exposure
- An intrusion exposed customer records and facial recognition training data.
- **Lessons Learned**: Securing biometric and proprietary model data requires encryption, third-party vetting, and strict access policies.

## Facebook AI Ad Targeting Vulnerability (2019)

- **Risk Type:** Inference-based privacy risk
- Research demonstrated that Facebook's ad system could infer sensitive user attributes without consent.
- **Lessons Learned:** Privacy-preserving methods such as differential privacy and federated learning should be applied to reduce inference risks.

## Lessons and Best Practices

Across these examples, several important principles emerge:

- Implement data lineage tracking and model version control to maintain visibility over inputs and outputs
- Monitor models post-deployment for drift, bias, and anomalies
- Apply granular access controls to prompt data, API calls, and intermediate outputs
- Leverage advanced encryption to protect sensitive data at rest, in transit, and in use to assure data sovereignty and reduce critical data exposure and/or manipulation
- Align AI data pipelines with change management and patching processes
- Review regulatory compliance at each stage of the AI system lifecycle

These incidents highlight how traditional controls are often insufficient when applied to AI systems. Organizations must adopt layered protections that consider the unique behaviors and risks of AI-based workflows.

# Harnessing the Future: Potential Challenges and Direction

## Technical Challenges in Securing AI Systems

Effectively protecting data in AI systems requires addressing foundational software and infrastructure issues. These include the continued use of end-of-life (EOL) operating systems, the absence of patch management practices, and the presence of unmitigated vulnerabilities. These baseline issues must be prioritized as part of a robust and proactive security strategy.

Encryption remains a critical layer of defense, yet many organizations rely on outdated or weak cryptographic methods. FIPS-validated encryption, endorsed by the National Institute of Standards and Technology (NIST), offers strong safeguards for sensitive information. Additionally, 256-bit encryption is widely accepted as a reliable standard for protecting data in transit (via SSL/TLS 1.3), at rest, and in use.

Data management challenges can also lead to limited visibility into data flows and security posture. Building a comprehensive data inventory, combined with proper classification and protection measures, helps ensure alignment with evolving regulatory obligations. A Data Protection Impact Assessment (DPIA) further supports the identification and mitigation of risks specific to AI data processing.

Access control must be reinforced with role-based access control (RBAC) systems that align user permissions with job responsibilities. ABAC Attribute Based Access Control for more granular access control for least privilege principle.

These controls are essential for limiting exposure and securing sensitive AI training, inference, and logging environments.

## AI Kill-Switch & Rollback

For critical sectors like finance, healthcare, and autonomous systems, AI must include a built-in "off switch." A kill-switch and rollback plan ensures that if a system goes wrong, it can be stopped or reversed instantly—protecting operations, customers, and compliance.

# Challenges AI Presents to Data Classification

Generative AI models and agents increasingly rely on large volumes of unstructured data. To manage this input, AI systems segment data into smaller chunks for more efficient processing. While this technique improves performance, it introduces four key challenges that disrupt traditional Data Loss Prevention (DLP) and classification strategies:

1. **Classification Bias at the Chunk Level:** AI classifiers often analyze data in fragments rather than in full context. This chunk-based approach may lead to inconsistent or inaccurate classifications, as the meaning and sensitivity of a segment can vary depending on the broader document or conversation it came from.

2. **Sensitivity to Chunk Size:** The accuracy of classification models depends on the amount of information available. Small segments may lack sufficient detail, resulting in misclassification. As the amount of data increases, classification becomes more precise and context-aware.

3.  **Mixing Chunks Within Context Windows:** When AI models process queries, multiple chunks from diverse documents may be retrieved and combined in the same context window. This aggregation can lead to unpredictable classification outcomes, as the resulting input no longer reflects the original document-level classification.

4.  **Duplication of Data Across Chunks:** Data reuse across emails, documents, and SaaS platforms is common. Boilerplate content, version control artifacts, and user copy-paste behavior contribute to high duplication rates. Studies from WAN optimization vendors suggest that more than 90 percent of enterprise data may be redundant at the chunk level, complicating classification accuracy and policy enforcement.

# Ethical Considerations and AI Bias in Data Security

AI-powered classification systems can unintentionally amplify bias and introduce ethical risk into security workflows. These systems often learn from historical data or external models that may not reflect organizational values or compliance standards.

Bias may result in overclassification or underclassification based on language, format, or user group—leading to unfair restrictions, privacy violations, or noncompliance. For example, AI trained on biased training sets may disproportionately flag certain data types as risky, or miss sensitive content in unfamiliar formats.

Further risk arises when inferred metadata or classification tags leak into downstream logs or outputs. AI-generated summaries or annotations may contain sensitive information not present in the original input. This semantic leakage undermines classification boundaries and can expose protected attributes to unauthorized users.

To mitigate these risks, organizations should:

- Regularly audit classification outputs across demographic, contextual, and content dimensions
- Incorporate bias detection and correction mechanisms into model pipelines
- Implement differential privacy or redaction techniques for AI-generated outputs
- Align AI classification behavior with legal requirements and ethical guidelines

# Future Trends in AI and Data Security

New technological paradigms such as quantum computing and federated learning are reshaping the future of AI security. While these innovations offer substantial benefits, they also introduce new risks and architectural challenges.

## Quantum Computing

Quantum computing threatens to break many current encryption schemes, particularly RSA and elliptic curve cryptography. As quantum capabilities mature, attackers may decrypt long-lived AI datasets, inference logs, or secure API traffic.

To prepare, organizations must adopt quantum-resistant cryptography such as lattice- or hash-based algorithms. Including, utilizing leading Security Standards such as NIST National Institute of Standards and Technology's.  Additionally, AI environments should update key management practices, enforce forward secrecy, and reclassify data assets based on quantum threat exposure.
New technological paradigms such as quantum computing and federated learning are reshaping the future of AI security. While these innovations offer substantial benefits, they also introduce new risks and architectural challenges.

## Federated Learning

Federated learning enables decentralized AI training across edge devices and organizations, without transferring raw data. While this model preserves privacy, it introduces new vectors for tampering, model poisoning, and gradient leakage.

Ex. Pilot Programs:
- Google GBoard
- Medical federated learning in EU projects

Securing federated learning pipelines requires:

- Secure aggregation techniques
- Differential privacy protections
- Cryptographic proofs to verify model integrity
- Strong contributor authentication and provenance tracking

## Multi-Modal AI Risks

Modern AI systems are increasingly working with multiple types of data—text, images, video, audio, and even sensor inputs. While this makes them incredibly powerful, it also introduces new risks. Insights from one data type can unintentionally reveal sensitive information from another, creating what's known as cross-modal leakage. Organizations should enforce clear standards and isolation controls to prevent such cross-modal leaks.

These developments demand that AI security architectures become more adaptive and forward-looking. Organizations should invest in post-quantum research, collaborate on federated AI standards, and embed resilience across cryptographic and infrastructure layers.

### Zero Knowledge Proofs (ZKPs) for AI Assurance

Zero knowledge proofs let organizations show that an AI model follows agreed rules—like avoiding certain datasets or applying differential privacy—without sharing sensitive data or revealing full model details. This approach gives regulators and partners confidence that models are compliant, while keeping proprietary information safe.

### AI Disaster Recovery / Continuity Planning

Organizations plan for system outages, but very few have playbooks for "what happens if your foundational model is poisoned or irreparably compromised?" Continuity plans should define fallback modes—such as reverting to older validated models, adopting human only review workflows, or rapidly onboarding third party foundation model providers—while maintaining compliance assurances.

# Conclusion

Artificial intelligence is reshaping how organizations must approach data security. While traditional frameworks are built on the principles of confidentiality, integrity, and availability, these foundations must now evolve to address the complexities introduced by generative AI, decentralized learning models, and real-time inference systems. AI introduces new risks that require organizations to rethink how data is collected, protected, and governed throughout its lifecycle.

This paper has outlined several critical threats unique to AI environments, including model inversion, prompt injection, data poisoning, and the unauthorized use of AI tools, often referred to as shadow AI. These risks highlight gaps in existing security models and underscore the need for AI-specific controls such as those proposed in the CSA AI Controls Matrix, particularly DSP-25 through DSP-28. These additions address emerging risks that are not yet adequately covered by conventional data protection frameworks.

To defend against these evolving threats, security architects and compliance leaders must deploy a combination of proactive and reactive strategies. This includes improving data classification, implementing strong encryption and data minimization practices, and using AI technologies for threat detection, anomaly monitoring, and behavioral analysis. Proper data labeling and visibility are essential. Unlabeled or misclassified data can bypass protections and introduce systemic vulnerabilities.

Another emerging priority is preventing the misuse of sensitive or proprietary client data in third-party AI model training. As regulatory scrutiny increases, organizations must adopt enforceable policies around data tagging, contractual safeguards, and auditing of external model usage. Gaining visibility into how external vendors handle training data is critical for maintaining compliance, protecting intellectual property, and building trust.

Looking ahead, effective AI security programs will need to integrate explainable systems, enforceable governance frameworks, and quantum-resilient protections. Research into secure federated learning, model validation, and privacy-preserving AI techniques will continue to be central to developing responsible and secure AI ecosystems.

AI threats evolve fast, and static security measures can quickly become outdated. To keep defenses effective, organizations should set up a regular peer review cycle—say, every 6 to 12 months—to reassess safeguards against new attack methods. This ongoing validation helps ensure security controls stay adaptive and reliable over time.

Trustworthy AI, depends on robust, lifecycle-aware data security. Without strong governance, continuous monitoring, and embedded ethical safeguards, the reliability and legitimacy of AI systems cannot be assured. Security and compliance programs must evolve alongside AI innovation to ensure both progress and protection.

# References

- **World Economic Forum.** (2024). Global Cybersecurity Outlook 2024 (GCO). https://www3.weforum.org/docs/WEF_Global_Cybersecurity_Outlook_2024.pdf
- **Association of Corporate Counsel.** IP Challenges in a Data-Fueled AI World. ACC Docket. https://docket.acc.com/ip-challenges-data-fueled-ai-world
- **Wiz Research.** (2025). DeepSeek Database Exposure – Database Encryption Risk. https://www.wiz.io/blog/wiz-research-uncovers-exposed-deepseek-database-leak
- **Cloud Security Alliance.** (2024). AI Organizational Responsibilities: AI Tools and Applications – Prompt Injection Defense.

https://cloudsecurityalliance.org/artifacts/ai-organizational-responsibilities-ai-tools-and-applications

- **European Union.** (2024). EU AI Act (Pending Final Adoption). Legislative framework for trustworthy AI.
- **General Data Protection Regulation (GDPR).** Regulation (EU) 2016/679. https://gdpr.eu/
- **California Consumer Privacy Act (CCPA).** California Civil Code §§ 1798.100–1798.199. https://oag.ca.gov/privacy/ccpa
- **Health Insurance Portability and Accountability Act (HIPAA).** Retrieved from: https://www.hhs.gov/hipaa
- **Brazil – General Data Protection Law (LGPD).** https://www.planalto.gov.br/ccivil_03/_ato2015-2018/2018/lei/l13709.htm
- **ISO/IEC 23894:2023.** Artificial Intelligence – Risk Management Framework.
- **NIST.** (2023). AI Risk Management Framework (AI RMF 1.0). https://www.nist.gov/itl/ai-risk-management-framework
- **Cloud Security Alliance.** (2024). Cloud Controls Matrix v4.0.13. https://cloudsecurityalliance.org/research/cloud-controls-matrix
- **Cloud Security Alliance.** (2025). Compliance Automation Revolution (CAR). https://cloudsecurityalliance.org/car
- **OWASP.** (2023). Top 10 for Large Language Model Applications. https://owasp.org/www-project-top-10-for-large-language-model-applications
- **MITRE.** Adversarial Threat Landscape for Artificial-Intelligence Systems (ATLAS). https://atlas.mitre.org/
- **OECD.** (2019). OECD Principles on Artificial Intelligence. https://oecd.ai/en/ai-principles
- **UNESCO.** (2021). Recommendation on the Ethics of Artificial Intelligence. https://unesdoc.unesco.org/ark:/48223/pf0000381137
- **ENISA.** (2023). Guidelines for Secure AI Development. https://www.enisa.europa.eu/publications/securing-machine-learning-algorithms
- **FATF.** (2021). Guidance on Digital Identity. Financial Action Task Force. https://www.fatf-gafi.org/en/publications/Fatfrecommendations/digital-identity-guidance.html
- **ISO/IEC 27090** (Draft). Information Security for Artificial Intelligence.
- **ISO/IEC 42001:2023.** Artificial Intelligence – Management System.
- **Facebook Ad Targeting Leak.** (2019). Analysis based on public disclosures and reporting on AI-based targeting inference vulnerability.
- **Clearview AI Breach.** (2020). Public reports on biometric data access and related lawsuits.
- **Google DeepMind Retention Incident.** (2022). GDPR-related data retention concerns in medical AI collaborations.
- **OpenAI API Chat History Exposure.** (2023). Based on OpenAI disclosures regarding a bug in their API session management.
- **NIST.** (2023). AI RMF Companion Playbook. https://www.nist.gov/itl/ai-risk-management-framework

- **TensorFlow Privacy.** Google AI. https://www.tensorflow.org/responsible_ai/privacy/guide
- **OpenDP.** Harvard University. https://github.com/opendp/opendp
- **Microsoft SEAL.** https://www.microsoft.com/en-us/research/project/microsoft-seal
- **OpenFHE Library.** https://github.com/openfheorg
- **Enveil.** Homomorphic Encryption Solutions. https://www.enveil.com/
- **Paperclip SAFE.** Searchable Encrypted Storage. https://paperclip.com/safe
- **PYSearch.** (Open source SSE library). https://github.com/ko1o/PYSearch
- **Duality Technologies.** Secure Multiparty Computation Platform. https://dualitytech.com/platform/duality-collaboration-hub/
- **MP-SPDZ.** Open-source SMPC framework. https://github.com/data61/MP-SPDZ
- **IMDRF.** (2023). AI Guidance for Medical Devices. International Medical Device Regulators Forum.
- **FINRA.** Artificial Intelligence in the Financial Sector. https://www.finra.org/rules-guidance/key-topics/artificial-intelligence
- **NAIC.** (2023). Model Law for Insurance AI. National Association of Insurance Commissioners. https://content.naic.org/model-laws
- **DIRF:(2025**)-https://cloudsecurityalliance.org/blog/2025/08/27/introducing-dirf-a-comprehensive-framework-for-protecting-digital-identities-in-agentic-ai-systems
- https://www.pointguardai.com/blog/appsoc-named-in-gartners-new-ai-security-platform-technical-guide
- **NIST AI Risk Management Framework**. https://www.nist.gov/itl/ai-risk-management-framework
- **NCSL AI 2024 Legislation.** https://www.ncsl.org/technology-and-communication/artificial-intelligence-2024-legislation
- **NIST Post-Quantum Cryptography Standardization Project.** https://csrc.nist.gov/projects/post-quantum-cryptography
- **CSA blog.** https://cloudsecurityalliance.org/blog/2024/04/01/un-ai-resolution-eu-ai-act-and-cloud-security-alliance-s-recent-efforts-draft-white-paper-on-ai-organizational-responsibility-for-core-security