

## Research Article

# Responsible AI Data Architecture: Embedding GDPR and PII Compliance into MLOps Pipelines at Enterprise Scale

Narendra Mangala  
Data Engineer Manager

**\*Corresponding Author**  
Narendra Mangala  
([manganarendra2@gmail.com](mailto:manganarendra2@gmail.com))

**Article History**  
Received: 27.02.2026  
Accepted: 20.03.2026  
Published: 30.03.2026

**Abstract:** Data retention, usage, and protection are often at odds with the speed, agility, and scale of machine learning operations (MLOps) systems, burdening organizations with risks, penalties, and security breaches. To resolve this contradiction in the context of enterprise-scale MLOps at a global financial services institution, a data architecture design is proposed that embeds General Data Protection Regulation (GDPR) and personally identifiable information (PII) compliance into the design and governance of data pipelines and data used in model training and inference. The recommendations are grounded in GDPR principles and requirements, as well as PII management best practices. Closely integrated with the data lifecycle and with clearly assigned accountability, these recommendations enable compliance mandates to become a natural by-product of MLOps, guardrails for ethical production AI, and ultimately a competitive differentiator. AI is increasingly used in critical social functions—in finance, healthcare, law enforcement, and the provision of essential services—often with significant real-time consequences for individuals. Misuse, malfunction, and bias in these systems can therefore have catastrophic consequences, and these concerns have spurred the introduction of regulations and proposed laws that demand adherence to foundational principles of transparency, equity, accountability, and social benefit. With the United Nations projecting that at least 3 billion people will be covered by AI regulations by 2024, it is essential that organizations adopt these principles consistently and at scale in the training and operation of machine-learned models.

**Keywords :** Artificial Intelligence; Data Pipelines; Data Discovery; Data Protection; Data Security; Data Governance; Enterprise Architecture; MLOPs; Engineering; People; Processes; Technology; Policy; Regulation; GDPR; Personally Identifiable Information; Compliance; Data Governance; Responsible; Privacy; Auditing; Notifications; PII; PCI; CCPA; PIPEDA; ISO; IEC; ML; Machine Learning; Enterprise MLOPs; Development; Operations; Implementation; Design; Management; Automation; Research; Modelling; Testing; Training; Auditability; Explainability; Visualisation; Monitoring; Risk Management; Assessment; Legal; Andisit; Requirements; Procedures; Guarantee; Cloud; Modelling; Standards; CIF; CII; CIA.

**Copyright @ 2026:** This is an open-access article distributed under the terms of the Creative Commons Attribution license which permits unrestricted use, distribution, and reproduction in any medium for non-commercial use (NonCommercial, or CC-BY-NC) provided the original author and source are credited.

## INTRODUCTION

Compliance with the General Data Protection Regulation (GDPR), as well as adherence to industry standards requiring protection of sensitive data, is a fundamental concern in federally and commercially controlled industries. It is also increasingly important for organizations operating in jurisdictions engaged in cross-border data sharing and cloud-transit environments. The necessity for protection of PII (personally identifiable information) is widely understood. Less broadly recognized is the ability to detect, tag, and protect PII throughout the data architecture, including structured, semi-structured, and unstructured types, as well as data in motion and data at rest. Least-privilege control throughout the enterprise data architecture is essential for enabling privacy by design, yet naïve implementations may hinder usability, innovation, and business efficiency.

An enterprise MLOps practice must be founded on both a solid architectural foundation and a data governance framework that accounts for compliance requirements throughout the data supply chain. Data governance artifacts, including a data

catalog, support integration of these compliance requirements into enterprise-scale data pipelines. Principled classification, labeling, and metadata management provide a comprehensive understanding of the data estate; successful regulation and standard compliance rely on completeness and accuracy of these foundational elements. Active intelligence throughout the data pipelines supports assurance of compliance with regulatory obligations such as those arising from the GDPR and similar standards.

### 1.1. Research design

Tenets and unifying criteria of regulatory compliance; supporting controls; governing artifacts and practices; key roles and responsibilities; and channeling of the PII and privacy-by-design concepts throughout the data lifecycle, pipelines, modeling processes, and risk management function. The proposed Responsible AI Data Architecture integrates these elements as a semantic model guiding compliance-supporting decisions and controls. MLOps at enterprise scale provides the underlying architectural and governance framework linking legal and ethical considerations for data handling with management of PII and fulfillment of data subjects’ rights. GDPR defines high-level principles for lawful processing of personal information, but its complex stipulations and peripheral dependencies require detailed and concrete implementation to avoid breach, expose organizations to regulatory penalties, and compromise stakeholder trust.

MLOps at enterprise scale satisfies the GDPR and PII requirements in the context of data pipelines. Common stages of an MLOps data lifecycle, the touchpoints for privacy, security, and governance controls; the actors and actions involved, and the artifacts produced form the basis for embedding GDPR compliance and PII protection within MLOps pipelines at enterprise scale. The ensemble of practical measures, instruments, and dependencies on supporting functions comprise GDPR Alignment in Data Pipelines, PII Discovery, Protection, and Redaction Techniques, Privacy by Design in Model Development, and Compliance Automation in CI/CD for MLOps.



**Fig 1: Responsible AI Data Architecture**

### 1.2. Background and Significance

The General Data Protection Regulation (GDPR), adopted by the European Union in 2016 and in force since 2018, sets out requirements for organizations handling personal data. These requirements apply to organizations in the European Economic Area and in other jurisdictions with data localization laws or GDPR-like frameworks, as well as to any organization outside the GDPR territory that processes personal data of people in the European Economic Area. Servers, sensors, and other devices hosted in the region may also be subject to GDPR provisions even if their operators are not. The GDPR mandates the provision of protective controls to the data subjects—the people whose data is being processed.

Enterprise-scale machine learning operations (MLOps) provide a comprehensive set of centralized components, artifacts, and mechanisms that manage, enable, and govern the machine learning (ML) lifecycle at scale. As MLOps cover the end-to-end ML lifecycle, and data pipelines are at its core, they should also embed the necessary controls to ensure compliance with relevant regulations, like the GDPR, and the protection of sensitive data types, like personally identifiable information (PII), within MLOps pipelines. Implementing GDPR requirements directly into the pipelines is critical to ensuring that data subject rights are being addressed—after all, it is the pipelines that operate on and provide the data for training models and for inference. Enabling and supporting a compliant, privacy-preserving approach for the development and deployment of ML models at enterprise scale is an essential contribution of MLOps indeed.

Equation 1: Consent basis

Let:

$C(P)$ = indicator that the data subject gave valid consent for processing  $P$

Then:

$$C(P) = \begin{cases} 1, & \text{if valid, specific, informed, and active consent exists} \\ 0, & \text{otherwise} \end{cases}$$

#### Derivation

The article says every pipeline activity involving personal data must have a lawful basis recorded in metadata. One possible basis is consent.

So for any pipeline step  $P$ , if the legal record says the person consented to that exact purpose, then the consent indicator is 1.

Interpretation in pipeline terms user agreed to use of their data, consent is recorded, purpose in metadata matches actual processing.

## 2. Theoretical Foundations of Responsible AI Data Architecture

Political and legal considerations exert a strong influence on the architecture of an organization's data environment, filtering the available design decisions. In the domain of data privacy, the general requirements of the European Union's General Data Protection Regulation (GDPR) shape the construction and operation of AI systems. These regulatory demands inform and motivate specific requirements and mechanisms for artificial-intelligence-backed data systems. The implementation of EU data protection legislation is no longer simply a box-ticking exercise for the MLOps team; EU security and privacy demands are now business enablers that define the breadth of AI capabilities across the entire organization.

Conversations with stakeholders from a number of regulated industries revealed four core tenets of principled data architecture: data minimization, purpose limitation, integrity, and accountability. Echoing the principles defined by the GDPR, these imperatives require that data should be collected and processed only in sufficient detail to satisfy the legal obligations of the business at that point in time. Any data beyond that minimum should be collected and stored only if absolutely necessary for the business, and all stakeholders should be aware of the nature and usage of the data throughout its life cycle. For relationships involving supra-national parties, effective safeguards need to exist to cover cross-border transmission; organisations must also be able to substantiate those safeguards when necessary.

### 2.1. Legal and Ethical Imperatives in Data Handling

GDPR principles, data minimization, purpose limitation, integrity, and accountability support responsible data handling. Legal and ethical considerations impose substantive duties on data-handling organizations. General Data Protection Regulation (GDPR) Article 5 codifies the core principles underlying processing activities: lawfulness, fairness, and transparency; purpose limitation; data minimization; accuracy; storage limitation; integrity and confidentiality; and accountability. GDPR envisions these principles embedded into data-processing designs and systems, not just documented as easy-to-replicate checklist-based narratives. Integration is understood to reflect responsible design: the easier an organization finds it to affirm compliance with a given requirement, the more likely it is that the requirement is being fulfilled in spirit. Data-minimization and purpose-limitation principles, together with a commitment to integrity and confidentiality, thus necessitate actionable architectures.

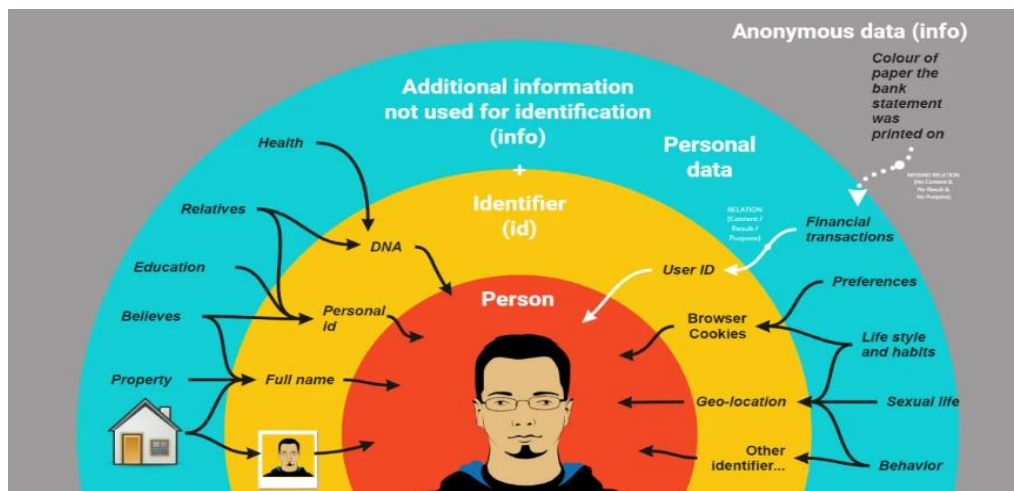
The principles delineate imperatives for responsibly handling data that support the arguments for a principled approach. In this context, GDPR compliance is understood as fulfilling the law's legal obligations, while in a broader setting, it encompasses also satisfying the more-general duty of responsible data handling.

### 2.2. Definitions: GDPR, PII, and Related Compliance Concepts

The General Data Protection Regulation (GDPR) governs the protection of personal data in the EU and EEA, and it regulates the export of personal data outside the EU and EEA areas. It recognizes that personal data protection is a fundamental right, and one of its objectives is to ensure that people's private life, family, home, and communications are protected against any interference. People have the right to freely express their views, and the EU institutions must respect this right. GDPR applies to the processing of personal data by a controller or a processor in the context of the activities of an establishment of a controller or processor in the EU, regardless of whether the processing takes place in the EU or not; to the processing of personal data of data subjects in the EU by a controller or a processor not established in the EU, where the processing activities are related to the offering of goods or services to such data subjects; and to the processing of personal data of data subjects in the EU who are located outside the EU by a controller or a processor not established in the EU, when the processing is related to the monitoring of their behavior.

Personally identifiable information (PII) is any data that can be used to identify a specific individual. Different organizations have different definitions of PII. For example, according to the National Institute of Standards and Technology, PII is any piece of information that can be used to distinguish or trace an individual's identity, such as name, social security number, biometric records, etc. The U.S. Census Bureau defines PII as any information about an individual that can be used to determine identity. The Joint Task Force released by the IEEE and the Internet Society states that PII is

a subcategory of sensitive information that can be used in online environments to differentiate users and provide access to services. Specific PII can also be divided into two broad categories: sensitive PII and nonsensitive PII. Sensitive PII requires extra protection (e.g., social security numbers, financial account numbers, and passport numbers), while nonsensitive PII (e.g., zip code and cultural heritage) can be made public without any consequence.



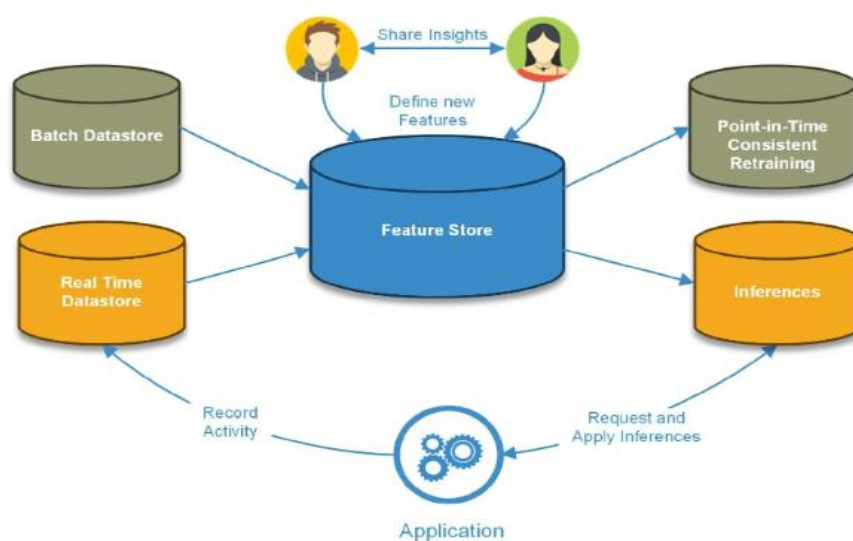
**Fig 1: GDPR, PII, and Related Compliance Concepts**

### 3. MLOps at Enterprise Scale: Architecture and Governance

Enterprise-scale Machine Learning Operations (MLOps) require architectures designed for frequent model development and deployment by multiple teams under a well-defined governance framework.

MLOps encompasses the full data and model lifecycle, from data discovery, ingestion, and preparation to model training, serving, monitoring, and retirement. The approaches adopted depend on the IT landscape and the organisation's privacy, security, and governance needs. Supporting frequent model development, automated retraining, and CI/CD pipelines can be challenging. Key aspects of enterprise-scale MLOps include architecture to accommodate multiple products or services, separation of duties for risk reduction, discovery, and cataloguing of available data along with automated controls to ensure privacy, data protection, and compliance with applicable regulations.

The data lifecycle in MLOps can be divided into stages: discovery, ingestion, preparation, training, serving, monitoring, and retirement. These stages are fundamental for understanding key touchpoints to implement privacy, security, and compliance controls during the data flow. Responsibilities for implementing and operating the controls across the lifecycle stages can be mapped to the standard roles of data governor, data steward, data custodian, privacy officer, and risk officer. An accountability matrix specifies the primary and secondary roles for each touchpoint, helping to segregate duties and assign appropriate level of oversight to risk-related privacy, security, and compliance controls.



**Fig 3: MLOps at Enterprise Scale**

#### 3.1. Data Lifecycle in MLOps

The data lifecycle in MLOps encompasses distinct stages from data creation or ingestion to utilization and retirement. Various controls ensure that these stages are managed according to the requisite governance, privacy, and security policies and procedures. The application and frequency of these controls may vary across the different data classes defined in the Cataloging and Classification section.

**\*Creation\*:** Data can be created internally, purchased from third parties, or sourced through public web crawls. Newly sourced or ingested datasets typically require a data management plan that defines their intended use and whether they will be made available externally. The data can then be captured in the Data Catalog, with accompanying logging, provenance tracking, and technical documentation for the ingestion process defining starting points along the Data Supply Chain. De-identification may also be applied at this stage, informed by the Data Protection Impact Assessment.

**\*Utilization\*:** Data in the Data Catalog may be used for training, testing, validation, or scoring of machine learning models under the approval and supervision of the Data Owner and Data Custodian. Usage requirements for sensitive data should be explicitly defined, possibly limiting usage to model development and requiring approval for actual model scoring in production. Model training pipelines, themselves part of the Data Supply Chain, should implement appropriate data minimization based on both the actual problem space and privacy considerations.

**\*Retirement\*:** Datasets that are no longer useful should be archived or deleted following the relevant Data Retention Policy. Archived datasets are preserved compressed and encrypted in a cost-efficient storage service. The Data Supply Chain tracks and identifies all copies of the retired dataset, including any derived datasets that have been released to external parties.  
Equation 2: Contract basis

Let:

$K(P)$  = indicator that processing  $P$  is necessary for performance of a contract

Then:

$$K(P) = \begin{cases} 1, & \text{if } P \text{ is necessary to fulfill a contract with the data subject} \\ 0, & \text{otherwise} \end{cases}$$

### Derivation

The paper emphasizes that each processing activity must have a documented lawful basis and purpose. If the operation exists to deliver a contracted service, then the contract basis applies.

Example using customer address to ship policy documents, using account data to execute a paid banking service.

### 3.2. roles and Responsibilities in Enterprise MLOps

MLOps consists of multiple stakeholders and diverse activities throughout the data and model lifecycle. A well-defined segregation of duties helps allocate responsibility and establish accountability across the MLOps construct. Infrastructure and platform teams are involved in the initial setup, as well as during support and maintenance. Governance frameworks and privacy enablement capabilities are defined and controlled by privacy and security functions. Data governance functions classify pipelines and designate the owners of used datasets. Data labeling and tagging services simplify the life of data scientists working on specific projects. Template-based deployment greatly reduces the risk of errors in production. Data engineers are responsible for deployment-ready pipeline implementations, model development teams for the models themselves, and support teams for their operational functioning.

An accountability matrix clarifies who is accountable for each MLOps stage. MLOps construct summaries show the handoff between teams. Data owners also retain ownership of specific pipelines. This model ensures a direct link between the dedicated topics of the MLOps governance framework and the involved MLOps pipelines, thus easing audit responsibilities. All these governance facets together provide the necessary mechanism to address touchdown areas across the lifecycle, such as privacy and security controls.

## DATA GOVERNANCE AND CATALOGING FOR COMPLIANCE

Effective compliance requires robust data governance rooted in strong data stewardship, as specified in the Data Governance Authority's annual report.

Governance can be defined as the establishment of policies, along with the continuous monitoring of their proper implementation, by the people responsible for a given area, so as to ensure compliance with regulations, goals, and objectives. Put differently, governance delineates Authority, Responsibility, and Accountability. In an enterprise context, data governance refers to the governing body, steering committee, and supporting organization charged with the custody, control, security, quality, availability, compliance, and KRI/KPI monitoring of an organization's data assets. Properly established authority, partnership, stewardship, and accountability roles and responsibilities are fundamental to effective

data governance in MLOps systems.

Data governance artifacts for compliance in MLOps systems serve as enablers and foundations for the integration of privacy, security, and regulatory requirements along the data lifecycle, from creation and aggregation to storage, use, sharing, and deletion. Cataloging practices supported by governance artifacts facilitate data discovery and accessibility and ideally feed into an organization's data catalog.

Data classification is the process of evaluating data assets to assign a classification label from a defined scheme—typically a restricted set of terms with a clear hierarchy—and, when appropriate, a sensitivity label that reflects special handling considerations. Automated tagging, powered by both steering-level classification input and natural language processing techniques, can streamline classification workflows.

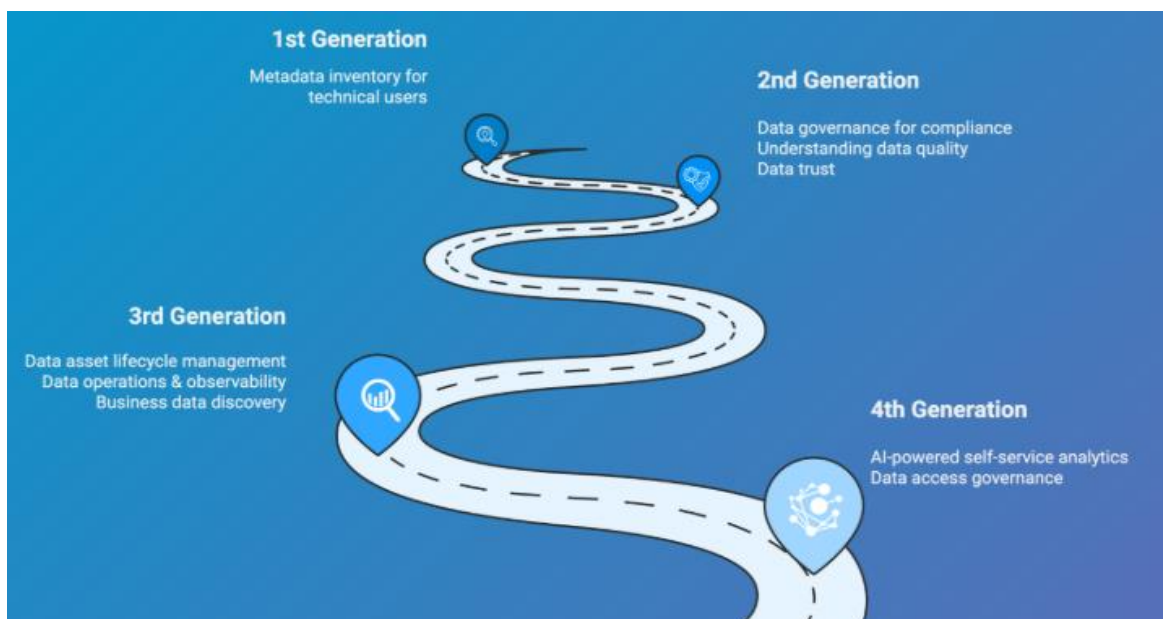


Fig 4: Evolution of Data Catalogs

#### 4.1. Data Classification and Labeling

A comprehensive classification scheme for enterprise data enables effective governance and control. A broad classification (e.g., Unrestricted, Internal Use Only, Restricted, Highly Restricted) provides an initial outline, while a more detailed sensitivity label taxonomy addresses regulatory and industry-specific considerations. Automated tagging during processing improves catalog comprehensiveness and reduces manual effort.

Data classification incorporates structured systems like the US Federal Information Processing Standards (FIPS), which define data as Unclassified, Sensitive (FIPS 199) or Non-Sensitive (based on impact levels). Labeling employs levels such as Low, Moderate, and High. Regulatory frameworks like GDPR, CCPA, PCI-DSS, HIPAA, and other applicable laws or regulations further refine classification with industry-specific considerations.

Machine-learning (ML)-powered computer-science models, Natural Language Processing (NLP) techniques, and ontologies also facilitate data classification. Although precise tagging remains challenging, output accuracy can be enhanced through a training corpus. Additionally, organizations' data-cleansing operations can identify source locations and assign preliminary classifications while updating the catalog. Further data classification occurs via lineage capture through data lineage tracking throughout pipelines. Proper column-name conventions can also provide semantic meaning. Tagging can be automated via custom-built, open-source, or commercial Classification (C) as a Service solutions using a combination of Natural Language Processing (NLP), Machine Learning (ML), knowledge graphs, and ontologies. Service selection requires careful evaluation of coverage, accuracy, performance, and maintenance costs. Automated metadata extraction can further reduce manual tagging effort but requires business-rule identification to establish patterns for the captured labels. Anonymization operations can also provide sensitive-context indications.

Equation 3: Legal obligation basis

Let:

$O(P)$  = indicator that processing  $P$  is necessary to comply with a legal obligation

Then:

$$O(P) = \begin{cases} 1, & \text{if } P \text{ is required by law or regulation} \\ 0, & \text{otherwise} \end{cases}$$

### Derivation

The article is about embedding compliance into enterprise pipelines. Some data processing happens not because the company wants it, but because law requires it, such as retention, audit, reporting, or fraud controls. That yields the legal-obligation indicator.

Example retention of transaction records for statutory periods, breach-related logging, regulatory reporting.

### 4.2. Metadata Management and Provenance

Well-defined machine learning model metadata is essential for effective governance and compliance. A formal schema defines required information for model management and data management integration. Provenance tracking at all stages gets implemented to create detailed data flow and transformation records.

Governance, compliance, and auditing standards mandate detailed process and data flow documentation. Whereas metadata typically addresses only model training artifacts, reproducibility requirements necessitate tracking records for all in-flight data and how models utilize that data. Captured metadata should ideally cover all data sources and transformations, listing origin, changes, and intermediary artifacts. Provenance information structures—essentially directed graphs—record such data processing details and support lineage-aware applications. Pipeline provenance implementation deploys a system that registers high-level data processing steps in a central repository automatically with appended details. The implementation specifics vary depending on the type of step and the code written for them. For example, an extract-transform-load step has different provenance registration code than a feature generation step.

The architecture further ensures that the pipeline user can append provenance information about datasets created as intermediaries in a step (like batch-job intermediate tables) or datasets that are overwritten or used only as targets. In the cloud context, terraform scripts and deployment jobs help register provenance of resource creations. and production deployments are updated to register that code and infrastructure changes happened in that platform-process step. For other steps where none of this is automatically possible, provenance appending still remains a manual process. These variables demonstrate high provenance completeness, making it suitable for GDPR compliance.

### 4.3. Data Retention, Minimization, and De-identification

Data retention policies determine the duration for which data remains accessible, establishing “minimum duration” and “maximum duration” for all data in the enterprise data catalog. Minimum duration is informed by detected risk exposure in the processed data, while maximum duration by data utility and value in the enterprise context. For time-series data, retention policies have to be defined with consideration for seasonal needs. Data subject requests for erasure are processed in a timely manner and escalated if involving sensitive data. Those treated as “third-party” data, including that under a data transfer agreement, are flagged and monitored.

Data minimization rules filter out unused columns in data at source, where applicable, with justification required for filtering at later stages. Obliterative de-identification is applied at source and at the time of exposing data for analytics by masking or removing PII and sensitive data detected by privacy-preserving detectors and other service-layer classification. Non-obliterative de-identification is also supported, such as through tokenization. Retention periods for de-identification tokens are set according to the tokens’ reconstruction risk, and cross-reference tokens for data matching and combining need to be maintained throughout valid relational query periods.

## 5. GDPR Alignment in Data Pipelines

Pipeline designs in a regulated environment must satisfy legal requirements for data processing. Select GDPR bases for processing, document pipeline purposes, and ensure that constraints on the use of data are enforced. Implement controls supporting the exercise of data subject rights (access, rectification, erasure) and automate steps that require manual involvement.

GDPR compliance revolves around six lawful bases for processing personal data, defined in Article 6. MLOps stakeholders can consult established privacy notices to assess whether standard contractual clauses for data shared with third parties meet the processing requirements. Operations that do not fall under any basis cannot lawfully take place. The regulation further stipulates that personal data must be processed for specified, explicit, and legitimate purposes. These must therefore be documented for each pipeline. Once defined, purpose limitation rules must be enforced. For example, data ingestion processes should block data that do not match the documented purposes. Similarly, personal data from databases can only be used if explicitly stated in the purpose description.

Data subject rights can also be mapped to technical controls, serving as monitors for the automated health check of

pipelines, detecting situations of non-compliance and triggering appropriate actions. Controllers must facilitate the exercise of such rights and complete requests within stated timeframes. Requests for access, rectification, and erasure that require execution by pipeline processes must therefore be defined. Any manual operation supporting these requests and its compliance with the GDPR timelines must be captured as part of the deployment procedures.

Equation 4: Vital interests basis

Let:

$V(P)$  = indicator that processing  $P$  is necessary to protect vital interests

Then:

$$V(P) = \begin{cases} 1, & \text{if } P \text{ is necessary to protect life or physical safety} \\ 0, & \text{otherwise} \end{cases}$$

### Derivation

Although less common in enterprise MLOps, GDPR includes this as a lawful basis. So if a processing event is required to protect a person's life or safety, then this basis is active.

Example emergency disclosure to prevent severe harm.

### 5.1. Lawful Basis for Processing and Purpose Limitation

Analysis of the processing activity within the data pipelines identifies the underlying lawful basis and precisely documents the purpose for each processing operation. The pipeline design and implementation pipelines represent the workings of both aspects.

The GDPR model recognizes the basic principle of "processing shall be lawful only if and to the extent that at least one of the following applies" and describes six possible bases for processing of personal data. These lawful bases for processing (Article 6 of the GDPR) must be identified, clearly stated and justified in the processing activity metadata record. In addition to Identification of the lawful basis for processing, Article 5(2) of the GDPR mandates that onus rests with the controller to demonstrate GDPR compliance, confirming that requirements described in the previous paragraph are insufficient.

Automation of all compliance-sensitive design aspects is an essential goal for pipelines. Incorporating an immutable policy administration point into pipelines can support fulfillment of the purpose limitation requirement, by restricting and controlling data destination capabilities based on a ruleset assembled from the aforementioned recorded purposes and the designated lawful basis for each downstream processing operation. Such a control could support a "data embargo" mechanism, preventing requests and/or measures that would result in transfer of data outside a region for which such data capture operation has not been authorized.

### 5.2. Data Subject Rights Implementation in Pipelines

Respecting the rights of data subjects imposes specific requirements on MLOps pipelines. All parties involved need to be informed when a request is received and the pipeline architecture must include the ability to respond in a timely fashion. Pipelines must be designed to provide the required information (for access or rectification requests), to effect data deletion and for the output to be delivered in a secure manner.

Moreover, changes made by data subjects must be properly reflected within the pipeline and consequently in all other dependent or consumer systems. Incorporating these requirements into standard operating procedures can help guarantee that they are met consistently.

Testing these flows with appropriate tools and at an acceptable frequency demonstrates implementation of the respective rights at a regulatory level.

### 5.3. Data Localization and Transfer Mechanisms

GDPR Article 44's requirement that any processing of personal data by an organization located in the EEA must occur within the EEA or in an "adequate" country implies that pipelines must be built on a foundation of data localization or, where feasible, compliant cross-border transfer mechanisms. Fully automated solutions differ from standard solutions: If data must be moved out of the EEA for processing, GDPR control processes must be followed; if data must be moved to an entity within the EEA but owned by an entity not covered by the GDPR, the controls must also be followed. Localized solutions can be considered "standard" solutions that an organization can use without additional governance layers. Although transferring the data can incur some overhead, a centralized policy governing data transfer to a service allowing more complex algorithms should very much decrease the extra layers of governance.

Some cloud service providers enable the customer to encrypt data before it arrives at rest on their service using encryption

keys that they do not control (the customers). In this way, the cloud provider cannot access the data in an intelligible format. If the container of the encryption keys and the customer key management service is outside of the provider location of the data, it should also not force a GDPR data transfer—even if the keys are being localized. The misplacement of keys should block access to the data, so the organization is automatically enforcing the idea of ownership and authority over the data.

## PII DISCOVERY, PROTECTION, AND REDACTION TECHNIQUES

Responsible data handling in machine learning pipelines requires explication of personally identifiable information (PII) and the deployment of appropriate protective measures. The detection and inventory of PII in development, testing, and production data are essential for subject-matter expertise in ensuring least-privilege access, that information is appropriately tokenized or encrypted, and that sensitive data is shielded in training and inference. Therefore, a PII detection capability—focusing on PII discovery across all business domains and its timely incorporation into the PII inventory supported by information governance practices—is enacted.

Effective methods are employed to discover PII in data at rest, data in transit, and data in use. These methods consider coverage of the various types of PII, the approaches employed (such as dependency-parsing or pattern-matching techniques), the accuracy of the detection techniques, and the methods for maintaining an up-to-date inventory of PII. Any organization relying on third parties to support machine-learning services must apply the principles of least privilege. An identity-and-access-management (IAM) model (or a suitable variant) is designed to ensure least-privilege access across the board. IAM incorporates mechanisms and processes to implement role-based access controls, regular reviews of access rights, and a “zero trust” principle for all non-human accounts.

Appropriate measures must be deployed to protect PII in data at rest, data in motion, and data in use. Encryption is a widely adopted technique supporting data protection but can introduce latency in time-critical systems. Therefore, it is important to ascertain when encryption is appropriate, and to assess whether other privacy-preserving techniques should be considered—such as tokenization or anonymization. These techniques naturally have different applicability, dependencies, and performance profiles, and may impose different requirements on the data retention policy.

Equation 5: Public task basis

Let:

$U(P)$  = indicator that processing  $P$  is necessary for a task carried out in the public interest or under official authority

Then:

$$U(P) = \begin{cases} 1, & \text{if } P \text{ supports an official/public task authorized by law} \\ 0, & \text{otherwise} \end{cases}$$

### Derivation

This basis applies mainly where the controller exercises public authority.

In the paper’s architecture language, it would still be recorded in the processing metadata exactly like other bases.

Example government or regulator-operated systems, officially mandated analytics in a public authority context.

### 6.1. PII Detection and Inventory

Numerous techniques exist for discovering PII in textual data, ranging from regular expressions to trained machine learning models. Regular expressions, reliant on syntactic patterns specific to individual data items, may be sufficient for certain structured data but are difficult to maintain for general-purpose detection. Machine learning-based systems can offer high coverage with greater accuracy, yet require more computational resources and continuity of annotated training data to retrain, validate, and test.

Whichever approach or combination of approaches is used, coverage and accuracy must both be complete and integrated into a production system. The resulting PII inventory should be treated as a living artifact by keeping it continually up to date during normal pipeline execution and regularly verifying its accuracy.

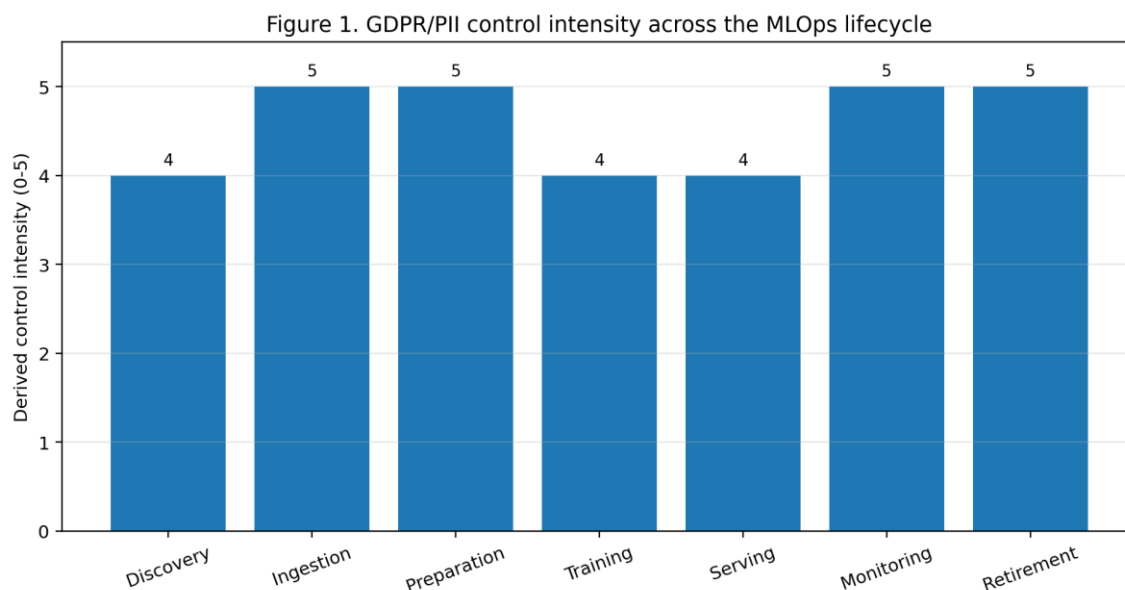
The goal in maintaining a PII inventory is coverage rather than 100% precision. The inventory must identify all PII in all data sources and active data at any stage of the pipeline. Coverage can be improved over time through root-cause analysis and rule addition of systematic false-positives. The effectiveness of rulesets such as those provided by AWS Macie or Azure Cognitive Services must also be monitored and improved.

### 6.2. Access Controls and Least Privilege

Access controls leverage Identity and Access Management (IAM) systems to limit data asset exposure in production environments. Access policies enforce the principle of least privilege—granting users the minimal level of access necessary to fulfill a well-defined business function. IAM builds role-based, resource-based, and attribute-based access control

models. Role-based access control maps predefined roles to a set of permissions, defining which actions a user can perform on designated resources. Resource-based access control facilitates self-service across a resource but requires an external access management function. Attribute-based access control takes contextual attributes into account to support dynamic conditions, enabling vacation override policies for role-based access control.

IAM privileges should be reviewed regularly by designated groups, typically the data asset owners, to align with the current needs of the business function. Privilege creep occur when users inadvertently accumulate access to multiple data assets across their lifecycle. Separation of duties is implemented by maintaining appropriate segregation of duties during the data lifecycle for actions such as production and acceptance of the data, implementation of new versions, and resetting of access controls. Wherever possible, IAM mechanisms prevent access to sensitive data within tools that support data life cycle operations such as data labeling and masking.



### 6.3. Encryption, Tokenization, and Anonymization Strategies

Four main strategies exist for protecting PII in data at rest: encryption, tokenization, anonymization, and pseudonymization. Encryption transforms data using a cipher and key into a ciphertext unreadable without the key. Data encrypted with the same key can be decrypted back to its original form. Tokenization replaces PII with a unique corresponding identifier—its token—generated by a tokenization vault or service. The `_mapping_` of original values to tokens is stored securely, such that any token can be re-identified with the original value if access to the vault or service is allowed. Encryption is generally reversible, whereas tokenization directly links sensitive data with a context-specific identifier creation, requiring no ciphertext decryption. Key management is a critical aspect of secret-based systems, whether being encryption, masking, or tokenization.

## 7. Privacy by Design in Model Development

Commitment to privacy should influence model development workflow design, affecting data selection, construction, and preparation. Several tactics result in privacy-preserving training pipelines. Techniques for minimizing the amount of sensitive data fed into ML models can prevent revealing individual-level private information through inferences. Accepting reduced performance in exchange for robust data minimization ensures compliance and builds trust in AI systems.

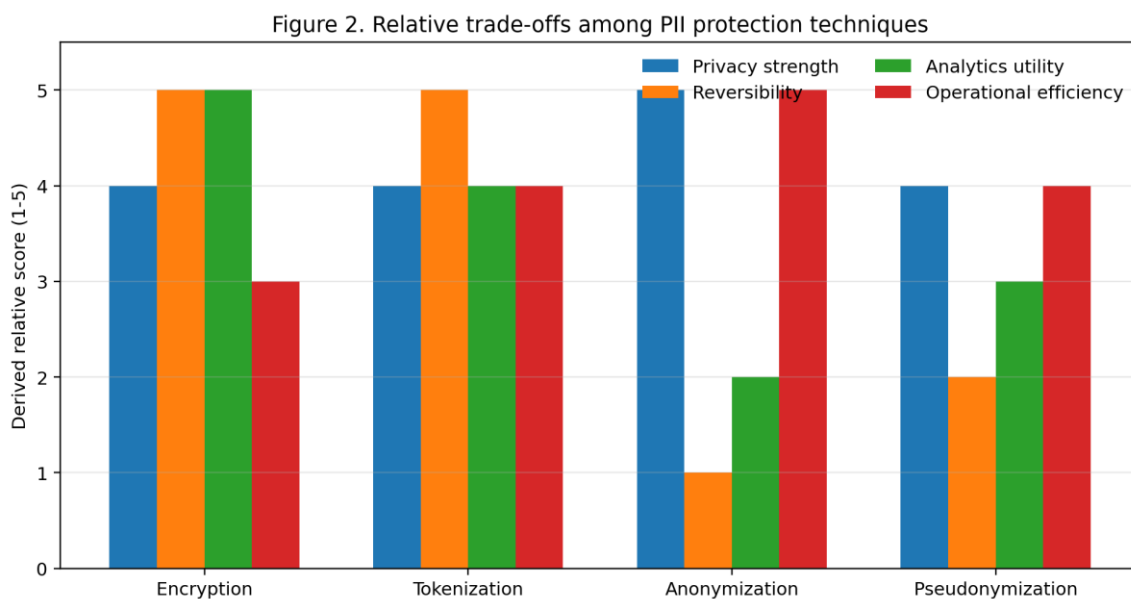
Increasingly common federated ML approaches involve training a global model while keeping data localized on individual data subjects' devices or data-holding entities' premises. Although federated learning is built on a collaborative model training paradigm that minimizes data concentration, it remains subject to federated learning privacy threats. Secure federation of non-collaborative data-holding parties can bolster privacy by design. It also allows pipeline architectures to scale in response to training data rarity, selection bias, and lack of generalization. Investment in a federated framework naturally integrates consideration of privacy threats into federated learning activities.

Privacy-preserving transformations should be considered at the feature engineering stage of ML pipelines, augmenting compliance-facilitating selection criteria grounded in ethics, business, and model performance. Selecting droppable features under differential privacy guarantees enables trustworthy privacy-aware deep learning; differentially-private prediction requires careful handling of sensitive features to mitigate risk of sensitive attribute inference.

### 7.1. Data Minimization in Training Pipelines

Training data is a prime candidate for minimization. Limits apply naturally to data sets for supervised learning and can also be applied in novel ways to semi-supervised learning models. To reduce the risk of leaking training data through adversarial attacks, it can help to minimize the number of training samples or their richness in sensitive information. Model developers can also explore different sources for proxy data that can stand in for the actual data to minimize risk during training. This requires applying such proxy data in tandem with a technique for validating whether or not the model is still performing in the desired manner. When a proxy is being used, trade-offs must be handled with care, particularly in sensitive contexts where prediction performance is highly correlated with privacy risk.

When replicating a model or service that has been in production, privacy risk levels may rarely be as strong of a concern when utilizing known data sources. Nonetheless, even in these situations, it is still good practice to minimize not only the volume of data being retained but also test predictions against an external source and minimize the raw amount of information being retained whether against data attacks or model inversion attacks.



### 7.2. Secure Federation and Federated Learning Considerations

Privacy-aware architectural designs and secure federated learning support model training with sensitive data while mitigating exposure. Secure federation for training is a technique that reduces data exposure by enabling training across multiple networks controlled by different organizations while keeping the member data within their respective network domains. Only the aggregated model information is shared with the federated learning coordinator. Hence, the training uses only metadata and model parameters of the models and not the individual input data that would normally be transmitted in the standard federated learning architecture.

In the classic federated learning architecture, edge nodes that are personally producing local models from sensitive data share them with a federated learning coordinator located at the cloud network. A potential security threat would exist since the personal local models may constitute important information for an adversary. It would be more secure if local models were never directly transmitted to the cloud network but rather aggregated there. In secure federation, the threat is addressed by encrypting and sharing only the model updates, while the local models remain secure inside the edges' protected space.

### 7.3. Privacy-Aware Feature Engineering

Data in an MLOps context often must relay personally identifiable information (PII) or business-relevant sensitive data (such as credit card details) in its raw or derivable form. Decisions regarding whether this information should be used in training and model inference represent common points of conflict between the modeling team and the data privacy office. Some decisions may be clear. For example, an image recognition system for an API that accesses images of children should rely on highly constrained training data, such as those aggregated in an independent and multisource secure federated setting with controlled contribution selection and participation. Conversely, other attempts to apply the same principle across different domains may not be viable, and engineers may need to use the same information in the models while deploying business-friendly proxies—essentially, lower-dimensional surrogates—based on feature selection methods designed—from a data privacy perspective—to maximize the minimization of information leakage from the original dataset.

During the prior analysis, the trade-offs obtained through models trained exclusively on a proxy remain within sample performance—that is, predictions made with minimal deviations and with good accuracy. The objective of designing and enforcing these principles is to broaden their application by encouraging their adoption when their performance impact during the MLOps is minimized.

## 8. Compliance Automation in CI/CD for MLOps

Integrating regulatory requirements into CI/CD DevOps processes is essential for compliance assurance in MLOps. Regulatory obligations must be expressed as formal policies and tested or enforced using dedicated gates in the CI/CD pipeline. Furthermore, auditability and explainability artifacts, such as reproducible reports, explainability records, and tamper-evident audit trails, establish accountability and traverse detection and response to non-compliance.

Compliance assurance requires integrating dedicated test suites into the CI/CD pipeline, covering the significant compliance controls in MLOps (attributes, techniques, and monitoring integrated in previous steps) and supported by formal regulatory frameworks. The results of these tests must generate a compliance report for the system. The test plans and results should be recorded to produce a reproducible audit artifact. Additionally, the responsibility matrix for MLOps systems must integrate compliance with regulatory obligations from the system definition and throughout its lifecycle. The CI/CD deliverables should include the confidentiality, integrity, and availability of the information processed by that system. These conditions should determine the formal requirements for the architecture that implements the system.

Compliance requirements are a critical and non-negotiable part of every system and business operation. To ensure that every aspect of system development and operation is aligned with compliance demands, the Policies-as-Code policy type should be defined and integrated within the CI/CD pipeline. These policies constitute business rules that determine what must happen throughout the execution of the pipeline—preventing the release of any artifact or the successful operation of any system that does not comply with the defined principles. These business rules provide guidance on the types of documentation generated and deliverables produced.

### 8.1. Policy-as-Code and Compliance Gates

Policy-as-code principles and reinforced regulatory compliance gates during continuous integration are vital capabilities for any enterprise-grade machine-learning platform. Firm policies regarding access control and data operations should be codified and enforced at different stages of the continuous integration/deployment pipeline.

Importantly, simple acceptance criterion documents and checklist-based compliance activities within the continuous deployment pipelines are insufficient for enterprise MLOps deployments. Manually executed privacy audits tend to be infrequent and prone to errors or oversight. Automated testing provides higher assurance that these regulatory compliance conditions hold for individual changes, yet manually created test suites can also drift over time. Guidelines can degrade into checkboxes, so any enterprise MLOps pipeline must enable active alerts or even gate a deployment if the policy is violated. Compliance auditing becomes meaningful only if everyone realizes that a violation will trigger the entire pipeline to fail, leading to timely correction before the affected change reaches production. Supporting these compliance gates within the continuous integration phase prevents manual review overhead and drift associated with individual test suites while ensuring that loss-averse production changes are made with the maximum permissible support. Multiple supporting technologies enable such automation within pipelines, each offering a mix of pros and cons that can be assessed during the design phase.

GDPR principle	Primary touchpoint	Representative control	Expected evidence
Lawfulness, fairness, transparency	Discovery / ingestion	Lawful basis registration, privacy notices, metadata	Processing record, catalog entry
Purpose limitation	Ingestion / serving	Policy checks restricting downstream use and transfer	Purpose policy, enforcement log
Data minimization	Preparation / training	Column filtering, proxy features, selective retention	Feature list, minimization rationale
Accuracy	Preparation / serving	Rectification workflows and data-quality checks	Correction ticket, test evidence

Storage limitation	Monitoring / retirement	Retention schedules, deletion/archival automation	Deletion log, archive manifest
Integrity & confidentiality	All lifecycle stages	IAM, encryption, tokenization, segregated processing	Access review, key policy
Accountability	CI/CD / monitoring	Policy-as-code, audit trails, reproducible reports	Compliance report, immutable log

**Table : maps the GDPR principles discussed in the paper to operational controls in enterprise MLOps.**

### 8.2. Automated Compliance Testing

Several subjects in the area of GDPR compliance need to be verified regularly, such as role-based access control or data retention. To establish confidence in ongoing alignment, a dedicated testing suite with the required coverage is created. Policy-as-code frameworks are useful here as well. The suite is typically incorporated into a CI/CD pipeline so that validation runs automatically after every code change and the results are reported in a clear, readable format. It is also good practice to review the test results periodically to gain a deeper understanding of the current state of compliance.

Certain data-related operations require the user to operate within a clearly defined legal framework. For instance, email correspondence must be archived in a way that ensures compliance with GDPR and other privacy-related regulations. When dealing with employees or customers in multiple locales, additional care needs to be taken. Non-public information and personally identifiable information should not be easily accessible to everyone, so a least-privilege access model should always be followed. Different forms of encryption, tokenization, or anonymization can also provide an extra layer of protection but have to be chosen carefully based on the use case.

### 8.3. Auditability and Explainability Artifacts

Privacy compliance requires a myriad of artifacts and documentation; many of them can be generated automatically. For regulatory alignment, for example, comprehensive audit logs are invaluable, providing a complete and accurate record of actions taken by data pipelines, including all accesses, changes, and deletions, together with all the data used and produced. When engineered correctly, these logs can be produced in a reproducible way at any point in time, when needed; a statement can be prepared that details how the requested data was processed and all accesses to the dataset in the specified period. These logs enable forensic analysis in the event of a privacy incident and can be leveraged for internal auditing and external regulatory scrutiny.

Similarly, explainability requirements for AI models mandate the production of interpretable records. In particular, tools for local model-agnostic interpretability such as LIME or SHAP can be applied on each inference request, associating a report to it. Privacy compliance for PII detection not only demands that detections be recorded, demonstrating evidence of operations on sensitive data, but also require that the records be interpretable. Reports generated during these operations satisfy these constraints and increase the compliance narrative by allowing greater transparency and insights on PII use—an area of increasing importance in the understanding of compliance posture.

Privacy compliance beyond artifacts can leverage these general principles: when an operation uses prescribed frameworks for auditing and explainability, the required outcomes are automatically fulfilled. Enabling efficient archiving, documentation, and decision-making increases the perceived effort under the three-layers of the effort currency, speeding future PII discoveries and articulation of details related to the classification of PII detection.

## 9. Monitoring, Audit, and Incident Response

Ensuring that compliance-related aspects of data pipelines are adequately monitored, that auditable records are maintained, and that the ability to respond to security incidents in a timely manner is also important from a Responsible AI perspective and serves as support pillars. Providing stakeholders with visibility into compliance coverage, applying in-depth monitoring for the aspects needing the greatest attention, and enabling rapid detection and response in incident situations mitigate the risks associated with these activities.

Ongoing assurance of compliance is vital for its effectiveness. A set of metrics dedicated to monitoring GDPR implementation—ideally expressed as a labeling scheme—should be defined. In addition to monitoring airplane inspections, monitoring data pipelines for GDPR compliance should employ anomaly detection techniques. High-value safety models should have monitoring dashboards that highlight, and facilitate remediation of, any functions that are not performing as designed. Monitoring dashboards should be curated for stakeholders with data privacy and protection responsibilities to provide them with coverage information about these models.

Comprehensive logging of the data pipeline system is required to achieve a favorable trade-off between the benefits and costs of explainability. The logs generated by the MLOps CI/CD environment should be built from a policy-as-code perspective, with a testing architecture able to trigger color-coded alerts for all detected compliance deviations in the short run.

GDPR compliance entails rapid detection and response to security incidents. An incident response playbook aligned with the GDPR should be established, covering notification timelines and essential analysis and remediation steps. For notifications, an anonymized record that addresses the demands of accurate and timely response should be maintained to satisfy the documentation requirements of article 33.

### **9.1. Ongoing Compliance Monitoring**

Ongoing assurance of compliance with regulations, internal policies, and contractual obligations is essential to managing legal, financial, and reputational risk effectively. In addition to the incorporation of compliance controls in pipeline design and operation, a program of continuous monitoring is key to detecting and addressing relevant changes. The identification of metrics for compliance-related KPIs is an integral part of the monitoring strategy: for example, data GIS classification discrepancies, policy antibody test results, minimum retention period breaches, and the ongoing review of access requests. These KPIs should be governed via a dashboard focused on the status of compliance controls and key risk indicators. Integration into a centralized governance platform provides the foundation for alerting and further anomaly detection.

By their nature, pipelines are systems subject to change. The addition of new sources or destinations may introduce previously unobserved risks, while the continual evolution of external factors such as regulation or business need requires an equally vigilant review of existing controls. For controls to remain effective, their implementation must keep pace with change. Failure to ensure this can result in poor outcomes either through over- or under-reaction, with potentially severe consequences. A set of change control principles specifically for compliance controls mitigates these risks. This list remains a living document, and each principle should be addressed proactively when evaluating a third-party amendment request or internal change request relating to a compliance control.

### **9.2. Logging and Tamper-Evident Audit Trails**

Sufficient logging is crucial for addressing not only the accountability obligations of GDPR Art. 5(2), but also those of Art. 24(1) that require appropriate technical and organizational measures to ensure compliance. Transparency caters to the principles of accountability and proportionality embedded in Art. 25, enables periodic reviews of compliance with data protection by design, and aids judges and other competent authorities in an independent investigation of potential infringements under Art. 57(1)(b). Incidents and breaches, including those pertaining to data security, must be documented as per Art. 33(5) and rec. 87; to achieve this, logging of incidents and forensics is essential.

The logs should capture all actions relevant to data retention, de-identification, segregation, monitoring, and third-party transfers. Log tampering must be prevented by creating a log structure that separates the log from the events themselves, ensuring logs are immutable and staff cannot disable logging functions, e.g., by using a third-party solution. In addition to being tamper-proof, logs should include integrity checks, such as cryptographic checksums on both individual events and in a Merkle tree for the entire log, allowing for public checks of integrity. Example checks for open-source solutions include public hash verification of chains and log removal detection using sparse Merkle trees.

### **9.3. GDPR Incident Response and Breach Notification**

A GDPR incident response playbook outlines actions to take in the event of a privacy incident regardless of timeline. It specifies near-term detection, reporting, investigation, notification, and decision-making responsibilities and procedures, followed by subsequent forensics activities based on the incident's scale. An additional annex defines precise timelines for specific data breach notifications related to the first and second phases. Playbook maintenance includes regular reviews and updates reflecting both GDPR requirements as well as changes to organizational policies and procedures.

GDPR notification requirements operate on a fixed timeline. They specify the window for notifying the relevant supervisory authority, the circumstances requiring disclosure to affected data subjects, and engagement of the data protection officer (DPO) when the breach is "likely to result in a high risk to the rights and freedoms of natural persons." Organizations make these determinations as part of their investigation. Decision-making responsibilities may also extend to legal counsel, public relations, and other functions.

## **10. Risk Management and Assurance**

Risk management and assurance touch all aspects of responsible data architecture as security, risk appetite, and maturity shape design decisions. Direct risk mitigation in data pipelines uses established assessment frameworks that produce concrete remediation actions. An additional lens ensures assurance-level risk factors associated with an organization's compliance posture are under control. For organizations integrating compliance within MLOps, third-party data use is especially important due to the inherent vulnerabilities that accompany external alliances.

Data pipeline risk assessments leverage one of the many available frameworks—NIST Special Publication 800-30, ISO 27005, FAIR, or others—and produce remediation actions. Natural candidates include the CIA triad, vulnerability-severity matrix, OWASP Top Ten, DREAD, and STRIDE. While these frameworks enhance understanding of pipeline security, identifying and managing other aspects affecting the compliance assurance level is equally crucial. Assurance concerns consist of anything that hinders or weakens the assurance provided to external stakeholders, such as lack of evidence, vulnerability to bribery or coercion, regulatory recognition, evidence reproducibility, and trust score.

### **10.1. Risk Assessment Frameworks for Data Pipelines**

Data pipelines are essential for delivering high-quality data to the many applications running on top of a data lake. Therefore, it is critical to ensure that all requirements in MLOps, data governance, and compliance with external legislation such as the GDPR are sufficiently met. A risk assessment framework allows a formal method for scoring the impact of non-conformity at the data pipeline level and for defining remediation plans. This approach can be leveraged not only in the context of compliance with GDPR but in any context where a formalized risk assessment is deemed useful for prioritizing actions or allocating budget.

The risk assessment framework employed follows the guidelines established in ISO/IEC 27005:2011 and ISO/IEC 27006:2015, as well as the risk-based approach described in ISO/IEC 27001. The applicability is further extended by MoSCoW prioritization, which incorporates the notion of urgency, alongside risk impact and likelihood. Thus, each assessment considers not only whether the requirement is met or not but also how severely it impacts the organization and how urgent its remediation is. Another dimensionality is added by Monte Carlo simulations, which allow estimation of the potential exposure of the organization per area. The identified risks and their corresponding remediation plans can be periodically reviewed and reprioritized as needed. Doing so can be incorporated into the existing Continuous Compliance Monitoring procedure.

### **10.2. Third-Party Compliance and Supply Chain Considerations**

Thorough due diligence on third-party data processors is vital to fulfill the GDPR requirement that processing takes place under the responsibility of the controller. Contractual clauses must transpose the controller obligations of Article 28 and supporting interoperability conditions of Article 46. However, a supply chain approach is also warranted, as third-party handling of personal data naturally extends beyond contractual relationships to potential partners, servant corporations, and any entity with access to the associated raw or processed data. The GDPR acknowledges the necessity for ongoing assurance in this area.

A risk-based framework tailored to third-party compliance is advisable, balancing the cost of maintaining the relationship with the vendor risk to the organization. Such factors as the burden of contractual compliance, the credentials available for direct due diligence, the potential reuse of data outside the profiling for data subjects (and the nature of such reuse), and the classification and sensitivity of data involved are all relevant. Consequently, while all vendors entail some degree of risk and require at least initial diligence, a subset must undergo more exhaustive baseline diligence, which may be satisfied through third-party report assessments (SOC, ISO 27001, BSI C5, etc.) endorsed by the organization.

Once initial diligence has been performed, compliance monitoring needs to happen on an ongoing basis, especially for vendors part of an outward-facing supply chain. Recent incidents affecting known providers of cloud computing services illustrate that security breaches within a single point may have global repercussions and that incidents are often detected and announced too late. As a result, continuous monitoring of third parties in a supply chain (however informal) is vital to maintain the required level of assurance over the data stored, processed, or transmitted.

## **11. Case Studies and Architectural Patterns**

Two case studies highlight practical implementations and lessons learned; architectural patterns distill common approaches for embedding GDPR and PII compliance into enterprise-scale MLOps pipelines.

**11.1. Case Study: GDPR-Driven Data Pipeline Redesign.** A major insurance company in the automotive sector sought GDPR alignment across its data assets and pipelines; sensitive information from motor insurance policies (dates of birth, addresses, and vehicle registration numbers) had been stored for analytics descopeing. An architectural review determined that residing in a non-compliant environment for months, even a few minutes, constituted a violation. The solution involved replacing traces of breach-specific data with anonymized equivalents. Synthetic identities generated by conference attendees' application relied on public datasets enriched with data from the European car registration database to preserve data-minimization and performance principles. Latency during upsert was impacted, and accuracy slightly reduced.

**11.2. Case Study: PII Shielding in Real-Time Inference.** Inference in a cloud-hosted chatbot for a multinational telecommunications operator required customer phone records and contracts; moderation was deemed insufficient to eliminate unintentional disclosure of PII. A "PII shield" was implemented to intercept and sanitize user input before routing to a personal AI model; shortcuts supporting general questions and event-specific requests could not afford latency

overhead. Entity and sentiment detectors automatically tagged PII; for detected credit card numbers stored in structured format, the last four digits, brand, and expiration date were unusable for chatbots and offered minimal risk for predictive models, allowing talk to facilitate engendered police and fraud reporting.

### 11.1. Case Study: GDPR-Driven Data Pipeline Redesign

A real-time data pipeline for Clickstream Analysis was retrofitted to ensure GDPR compliance, in particular Lawful Basis for Processing, Data Subject Rights, and Data Minimization. The previous design stored detailed user interactions for training click-through-rate prediction models and for generating a data mart supporting marketing initiatives. The marketing team was a secondary data consumer, while the models were not yet production grade. These factors weakened the justification for retention. Using live events stored only for reduced latency and immediate marketing impact was deemed acceptable. The architecture was modified to offload GDPR-sensitive attributes into an auxiliary store, restricting access and with additional protection applied on-demand. Latency contrasted with access controls, and shielding during real-time inferences was studied.

Data minimization is a key GDPR principle that requires limiting data collection and retention to the minimum necessary to fulfill a legitimate purpose. The real-time pipeline had initially stored complete clickstream events. These events supported quick marketing campaigns but were not essential for longer-term analysis. Retaining all events with GDPR-sensitive attributes was difficult to justify, especially since federated learning scenarios were being explored. As the click-through-rate prediction model was not production ready, the marketing team provided the strongest second use case. To ease management and augment on-demand protection, a complex document-store setup suitably shielded GDPR-sensitive attributes while minimizing additional latency. Actual shielding times and accuracy impacts were analyzed as an aid towards operationalization.

### 11.2. Case Study: PII Shielding in Real-Time Inference

Compliance imperatives that require PII protection may usefully be domestically mapped to privacy-in-ML machine learning architectures or PII-shielding techniques. At runtime, when decisions rely on training-anonymized models operating over PII-sensitive features, privacy-shielded feature values are fetched from a PII Protector readthrough and passed to the model as needed.

Pi-protection techniques guard against PII exposure by ensuring that no PII data are used in training the model while enabling model-based access to PII data at inference time. Keeping training data free of PII allows latent-space anonymization to be effective. However, PII-proofed features must still be supplied at runtime when the predictions require them. A runtime pass-through copy of the data minus a given set of PII-sensitive features is thus maintained, effectively performed by a PII Protector component gracefully shielding the PII channels.

## CONCLUSION

Responsible AI Data Architecture realizes its objectives by embedding GDPR and PII compliance patterns into the MLOps pipeline data architecture. Supportive data governance practices enable privacy-by-design integration within and alongside enterprise-scale pipelines. These contributions elevate MLOps beyond mere monitoring of regulatory requirements; compliance becomes a fundamental nonfunctional property at all stages of pipeline operation and evolution. A better-structured software architecture, based on principled design, offers greater operational resilience, clarity, and maintainability.

While the research focus is on fulfilling GDPR requirements, the mechanisms and procedures defined here also create a foundation for implementation of privacy frameworks in other jurisdictions. Others working on PII-free model training can draw on the GDPR-compliance patterns to enhance resilience and clarity in design, development, and operations. Broader adoption of the proposed principles should result in improved user trust in AI systems and bolster their deployment.

## REFERENCES

1. Pamisetty, A., et al. "Explainable AI Systems for Credit Scoring and Loan Risk Assessment in Digital Banking Platforms." *2025 IEEE 13th International Conference on Intelligent Data Acquisition and Advanced Computing Systems: Technology and Applications (IDAACS)*, IEEE, 2025, pp. 1478–1483.
2. Backes, C., R. Müller, and M. Schneider. "The Right to Be Forgotten: A Modern Approach to Data Deletion." 2019, pp. 177–210.
3. Mangalampalli, B. M. *Architecting Smart Health Economies: Data Fusion, Cognitive Automation, and Payment Integrity*. Deep Science Publishing, 2026.
4. Bathini, T., et al. "Responsible AI: Unpacking GDPR Obligations." 2020.
5. Nagubandi, A. R. "Future Directions in Autonomous Market Infrastructure." *Cognitive Financial Infrastructure: Designing Adaptive, Integrated Market Systems*, Deep Science Publishing, 2026.
6. Borodin, Y., et al. "Privacy-Policy-Aware Data-Centric Applications." 2021.

7. Garapati, R. S., et al. "The Evolution of Digital Payments: A Study on AI-Powered Transaction Monitoring Systems." *2025 3rd International Conference on IoT, Communication and Automation Technology (ICICAT)*, IEEE, 2025, pp. 1–8.
8. Camara, A., et al. "A Comparison of Data Auditing Algorithms for Private Cloud Services." 2021.
9. Kolla, T. "AI-Powered Data Catalog Systems for Healthcare Data Discovery and Governance." *South Eastern European Journal of Public Health*, 2024, pp. 2296–2311.
10. Chatzikokolakis, K., et al. "An Algebra for Privacy with Applications." 2015.
11. Gupta, D. K., et al. "Semantic Feature Learning Using Transformer-Based Deep Neural Networks." *2025 IEEE 5th International Conference on ICT in Business Industry & Government (ICTBIG)*, IEEE, 2025.
12. Dhanesh, N., et al. "Predictive Modelling with Privacy Preservation Measure." 2019.
13. Bandi, V. D. V. K. "Self-Optimizing Data Pipelines Using Machine Learning for Cloud Workloads." *Journal of Information Systems Engineering and Management*, vol. 10, 2025, pp. 1618–1636.
14. Nellutla, N. *MLOps 2.0: A Reference Architecture for CI/CD with Continuous Data Validation*. 2026.
15. Bargavi, N., et al. "Safeguarding Consumer Data in Digital Insurance: Legal Frameworks and Ethical Imperatives." *International Insurance Law Review*, vol. 34, no. S1, 2026, pp. 272–284.
16. Jeon, I., et al. "UnPII: Unlearning Personally Identifiable Information with Quantifiable Exposure Risk." 2026.
17. Kolla, S. H. "Future Directions in Trusted and Self-Regulating Enterprise Intelligence Systems." *Secure and Governed Enterprise Intelligence Platforms*, Deep Science Publishing, 2026.
18. Pahune, S., et al. "The Importance of AI Data Governance in Large Language Models." 2025.
19. Segireddy, A. R. *Cloud-Scale Intelligence for Financial Platforms: Adaptive Systems and Operational Artificial Intelligence*. Deep Science Publishing, 2026.
20. Asthana, S., et al. "Adaptive PII Mitigation Framework for Large Language Models." 2025.
21. Pallapu, S. R., et al. "GAN-Augmented Transformer Framework for Cross-Domain Video Style Transfer." *2025 International Conference on Communication, Computer, and Information Technology (IC3IT)*, IEEE, 2025.
22. Zhang, L., et al. "A BERT-Based Empirical Study of GDPR Compliance in Privacy Policies." 2024.
23. Kolla, S. H., et al. "Secure RAG Architectures with Small Language Models for Governance-Aligned LLM Deployment." *International Journal of Economic Practices and Theories*, 2026, pp. 166–179.
24. Mangalampalli, B. M. "Future Horizons in Sustainable and Adaptive Health Enterprises." *Architecting Smart Health Economies*, Deep Science Publishing, 2026.
25. Nagubandi, A. R. "Governance, Transparency, and Trust in Intelligent Financial Systems." *Cognitive Financial Infrastructure*, Deep Science Publishing, 2026.
26. Amershi, S., et al. "Software Engineering for Machine Learning: A Case Study." 2019.
27. Kolla, S. H. "Scaling Autonomous Execution Across Enterprise Services." *Secure and Governed Enterprise Intelligence Platforms*, 2026.
28. Zaharia, M., et al. "Accelerating ML Lifecycle with MLflow." 2018.
29. Krishnan, M., et al. "Human-in-the-Loop Hybrid Neuro-Symbolic AI Model for Reliable Data Engineering." *2025 IEEE Global Conference on Wireless Computing and Networking*, IEEE, 2025.
30. Kim, M., et al. "Data-Centric AI: A Systematic Review." 2021.
31. Bandi, V. D. V. K. *Autonomous Data Platforms: Converging AI, MLOps, and Cloud Engineering*.
32. Hummer, W., et al. "Benchmarking MLOps Platforms." 2019.
33. Sheelam, G. K. "Agentic AI in 6G: Revolutionizing Intelligent Wireless Systems." *Advances in Consumer Research*, 2025.
34. ISO/IEC. *ISO/IEC 27001: Information Security Management Systems*. 2022.
35. Surendra Yandamuri, U. "Governance, Security, and Responsible System Design." *Operational Intelligence Engineering*, 2026.
36. NIST. *AI Risk Management Framework (AI RMF 1.0)*. 2023.
37. Devayani, G., and K. C. Nagabhyru. "Wireless Sensor Networks and Digital Twins for Real-Time City Simulation." SSRN, 2026.
38. Voigt, P., and A. von dem Bussche. *The EU General Data Protection Regulation (GDPR)*. 2017.
39. Cate, F. H. "The Limits of Notice and Choice." 2010.
40. Tene, O., and J. Polonetsky. "Privacy in the Age of Big Data." 2012.
41. Abadi, M., et al. "Deep Learning with Differential Privacy." 2016.
42. Kairouz, P., et al. "Advances and Open Problems in Federated Learning." 2021.
43. McMahan, B., et al. "Communication-Efficient Learning of Deep Networks." 2017.
44. Barocas, S., et al. *Fairness and Machine Learning*. 2019.
45. Gebru, T., et al. "Datasheets for Datasets." 2021.
46. Raji, I. D., et al. "Closing the AI Accountability Gap." 2020.
47. Selbst, A. D., et al. "Fairness and Abstraction in Sociotechnical Systems." 2019.
48. Jobin, A., et al. "The Global Landscape of AI Ethics Guidelines." 2019.
49. Machanavajjhala, A., et al. "l-Diversity: Privacy Beyond k-Anonymity." 2007.
50. Fung, B. C. M., et al. "Privacy-Preserving Data Publishing." 2010.

51. El Emam, K., and F. K. Dankar. "Protecting Privacy Using k-Anonymity." 2008.
52. Aggarwal, C. C., and P. S. Yu. *Privacy-Preserving Data Mining*. 2008.
53. Lindell, Y., and B. Pinkas. "Secure Multiparty Computation." 2009.
54. Redman, T. C. *Data Driven: Profiting from Your Most Important Business Asset*. 2013.
55. Ebert, C., et al. "DevOps for AI Systems." 2021.
56. Chen, T., et al. "Data Validation for ML Pipelines." 2021.
57. Schelter, S., et al. "Automating Large-Scale Data Quality Verification." 2018.
58. Kandel, S., et al. "Wrangler: Interactive Data Cleaning." 2011.
59. Rekatsinas, T., et al. "HoloClean: Holistic Data Repairs." 2017.
60. Khayyat, Z., et al. "Big Data Cleaning." 2015.
61. Abedjan, Z., et al. "Detecting Data Errors: A Survey." 2016.
62. Kumar, A., et al. "Feature Stores for Machine Learning Systems." 2020.