



Automating Data Governance and PII Compliance Using Unity Catalog in AI-Driven Data Ecosystems

Ganesh Pambala

Independent Researcher, India

ganeshpambal@gmail.com

ABSTRACT: Data governance encompasses defining the roles, responsibilities, and accountability needed to safeguard data assets; enabling access control and usage monitoring; and supporting policy-driven data discovery, classification, protection, and retention. The increasing pervasiveness of Artificial Intelligence (AI) in analytics, data science, and machine-learning workloads highlights the importance of managing sensitive information and complying with legal requirements regarding PII. Numerous regulations compel organizations to prevent PII breaches while allowing analytics. Third-party cloud data platforms simplify AI-driven data ecosystems but pose a risk of PII exposure when sensitive information is shared across multiple environments, including untrusted external entities. These issues can be addressed through automated data governance, using policy-driven workflows that define and enforce policies related to PII and data governance.

Unity Catalog extends the data platform's capability to manage metadata across multiple cloud-object-storage accounts, implementing policy-driven automation for PII compliance and data governance through two approaches. The first approach automates key management and PII mapping to Data Loss Prevention tags, independent of an Identity and Access Management cloud service. The second approach enforces policies defined in an external Identity and Access Management service, with the cloud data platform as a service consumer rather than an IAM vendor. Implementation details gleaned from an enterprise production environment illustrate how Unity Catalog can automate PII-data governance and PII-compliance workflows.

KEYWORDS: Unity Catalog, Data Governance Automation, PII Compliance, AI-Driven Data Ecosystems, Data Access Control, Metadata Management, Data Lineage Tracking, Privacy Regulations (GDPR, CCPA), Sensitive Data Classification, Role-Based Access Control (RBAC), Data Security Frameworks, Automated Policy Enforcement, Data Auditing & Monitoring, Lakehouse Governance, Fine-Grained Data Permissions.

I. INTRODUCTION

The growing demand for AI-driven data products in customer experience, marketing, and product development has resulted in a significant expansion in the size, number, complexity, and maturity of data science and analytics ecosystems. Data assets, including sensitive personal identifiable information (PII), are becoming essential for new data-driven product differentiation and revenue sources. However, complex PII classification, identification, de-identification, reconciliation, and policy enforcement remain major reprioritized data governance challenges for organizations supporting predictive analytics and AI solutions.

Automating data governance and compliance using Unity Catalog can help address these challenges. The policies defined in Unity Catalog augment existing data discovery and metadata management capabilities to identify PII data types and enforce privacy policies. Data discovered through PII Rules configured in Unity Catalog can then support role-specific automation workflows, integrating identity and access management procedures for automated approvals and notifications related to data masking and access exemptions. A recent implementation illustrates how these capabilities enable policy-driven automation for sensitive data use in MI and how automation significantly reduces effort, risk, and time taken to satisfy such requests.

1.1. Background and Significance

To harness the potential of AI, organizations have started to build a modern, integrated data stack that is often distributed across multiple clouds and geographies. Such a data landscape enables them to extract insights from



machine learning models trained on a diverse mix of data sets. However, given the sheer number of users across different lines of business, organizations also face mounting pressure to ensure adequate data governance for its proper use. Data-related policies fulfil a key role in driving compliance and improving trust in the data ecosystem. Moreover, with the advent of Large Language Models, the risk of exposing personal data to AI models has only become higher. An evidence-based framework can be developed to create and manage Policies related to data discovery and classification, access control, compliance, auditing, and exception management—making these processes automated and consistent across the enterprise so that they do not remain a bottleneck for business.

Polices for the identification and cataloging of PII (Personal Identifiable Information), as well as for data de-identification and anonymization, can also be automated based on information from the Data Catalog. Finally, the enforcement of these policies can be integrated with an organization’s Identity and Access Management (IAM) tools to ensure that data custodians have clear visibility of governing roles and responsibilities. Using Unity Catalog and integrating its functionality with Simulated IAM Policy Functionality serve to demonstrate the approach and its potential—from design to automation to implementation and testing—along with the discussion of opportunities and caveats.

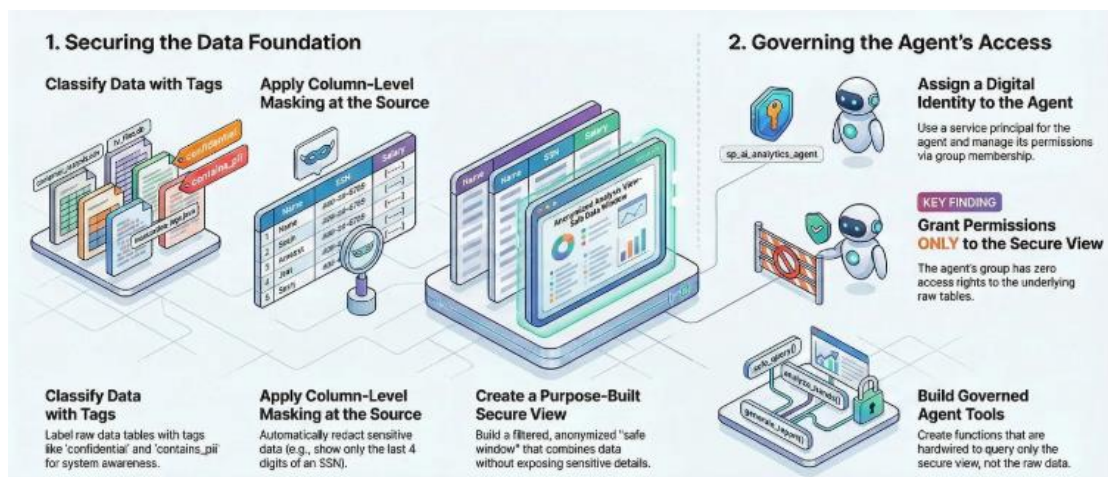


Fig 1: Building AI Data Governance with DataBricks

1.2. Research design

Evidence-based design complements the integrated discovery and interactive examination of policy domains. Privacy policy is governed persistently during the execution of AI operations, with an emphasis on automation and preventive controls based on policy-driven workflows. Administration is embedded in identity and access management systems for control and audit. Policy language supports expressions of data discovery, the applicability of privacy policy, and operational workflows for implementation, enforcement, deviation, and exception. Detection, classification, and cataloging of personally identifiable information leverage the classification system of data lakes and the data discovery capabilities of data warehouses. Support for de-identification and anonymization techniques is defined. Automation is achieved by endpoint preparation, code generation, other policy-driven automation opportunities, and repeatable tasks embedded in specialized data flows. Implementation in Unity Catalog enables wide applicability in demanding AI-driven data lakes while supporting operational stakeholders.

Equation 1: PII risk scoring equation

Step 1: Define the factors

Let

- R = re-identification risk
- E = identity exposure
- C = consent sensitivity
- A = auxiliary-information risk

Step 2: Assign weights

Suppose the importance weights are

- w_R, w_E, w_C, w_A



Step 3: Form the total risk score

A standard linear aggregation is

$$S_{PII} = w_R R + w_E E + w_C C + w_A A$$

Step 4: Optional normalization

If each factor is scored from 0 to 1 and the weights sum to 1,

$$w_R + w_E + w_C + w_A = 1$$

then $S_{PII} \in [0,1]$.

Step 5: Classification rule

Use thresholds:

$$\text{Class}(x) = \begin{cases} \text{Low sensitivity,} & S_{PII} < \tau_1 \\ \text{Medium sensitivity,} & \tau_1 \leq S_{PII} < \tau_2 \\ \text{High sensitivity,} & S_{PII} \geq \tau_2 \end{cases}$$

II. CONCEPTUAL FOUNDATIONS OF DATA GOVERNANCE

The Information Technology Governance Institute defines information technology governance as "the processes that ensure the effective and efficient use of IT in enabling an organization to achieve its goals." Information Technology Governance Institute (2011) Consequently, the goal of information technology governance is to create a framework that integrates all information technology domains—hardware, software, data and information, procedures and operations, personnel, finance, planning, and project management—into a unified management approach that addresses the organization's unique information technology needs. Therefore, it is possible to think of data governance as the part of information technology governance program that focuses on the data and information domain.

Data governance strives to meet the requirements of the data stakeholders and data accountable roles by establishing the necessary decision-making forums (i.e., data committees) that enable stakeholders' access to relevant, trustworthy data and information in a timely manner. Successful data governance relies on appropriate levels of decision-making authority and responsibility within the organization. These are addressed by defining the roles, responsibilities, and accountability for each of the process areas of data management—data discovery and metadata management, data quality, data integration, data warehousing and business intelligence, master data management, data security and privacy, data architecture, data operations, data support, and external data access. The design and operations of these processes are organized and executed through a clear delineation of sponsor, data owner, data creator, data steward, data user, and external data user roles.

2.1. Definitions and Scope

Data governance refers to the control, management, and oversight of data throughout its lifecycle, including the terms, conditions, and policies underlying its use. Compliance with applicable laws and regulations, such as those governing personally identifiable information (PII), is of primary concern — not only to avoid potential litigation for failing to sufficiently protect sensitive data but also to comply with expanding legal frameworks. As organizations become more reliant on data as a key enabler for creating business value, establishing an adaptive governance framework satisfying business needs becomes a priority. Such a framework defines roles, responsibilities, and accountability for the management of data and its utilization throughout the organization and aligns its policy objectives and business strategy with the evolving external regulatory environment.

Data governance shifts away from an initial focus on IT-centric stewardship or control of data toward a collaborative model in which the business takes ownership of the data and relies on IT for support. Essential elements in this shift are policy-driven automation and integration with an organization's identity and access management (IAM) systems. Applying the approach to processing sensitive business data enables compliance with regulations governing PII. Unity Catalog from Databricks serves as the underlying foundation, enabling organizations to manage enterprise data in a centralized manner, pivoting from an enterprise data lakehouse to a data marketplace to an AI-enabled enterprise.

Equation 2: Generalization equation

Step 1: Let the exact value be x

For age, $x = 27$.

Step 2: Choose a bin width

Suppose we generalize by decades, so bin width $b = 10$.



Step 3: Compute lower and upper limits

For interval-based generalization,

$$L = b \left\lfloor \frac{x}{b} \right\rfloor$$
$$U = L + b$$

Step 4: Substitute $x = 27$, $b = 10$

$$L = 10 \left\lfloor \frac{27}{10} \right\rfloor = 10(2) = 20$$
$$U = 20 + 10 = 30$$

Step 5: Write generalized output

$$G(x) = [L, U] = [20, 30]$$

Final derived equation

$$G(x) = \left[b \left\lfloor \frac{x}{b} \right\rfloor, b \left\lfloor \frac{x}{b} \right\rfloor + b \right]$$

For $x = 27$, $b = 10$, this gives $[20, 30]$.

2.2. PII and Compliance Paradigms

Personal data is defined as a 'particular type of information that is assigned to any individual and can be used to identify that individual and to distinguish them from others', while personal identifiable information refers to 'any information that can be used to identify an individual.' The definitions themselves vary depending on the jurisdiction or context considered. For example, the data protection regulation of a country/region consortium (such as GDPR in Europe) defines privacy explicitly. However, countries or regions that have not yet implemented their own regulation follow the PII listing published by the United States Department of Commerce when determining which set of data requires special protection. There is common use of solutions and tools (these can also be developed in-house) to identify and manage PII in enterprise data across organizations – especially, since in most national legislations, failure to comply can lead to costly penalties, as well as an organization damaging its brand image and reputation.

Data privacy compliance typically includes not only the data privacy regulation of a consortium of countries, but also those regions that have enacted their own data security regulation. The classification of data usage can be further split into production or non-production usage. The production usage ruling is strictest, and furthermore, violations in production usage can lead to very severe penalties. The non-production usage ruling, while also strict, offers a different means of protection through a single point of responsibility for PII usage across different areas of the enterprise. Additionally, the distribution of the security responsibility to the actual users of production data is supported through the permanent delegation of the security responsibility.



Fig 2: PII and Compliance Paradigms

III. RESEARCH SUMMARY

In AI-driven data ecosystems, data discovery, access control, and metadata management must be automated. Expensive manual processes—such as identifying which fields contain PII, flagging appropriate IAM roles, and actually applying IAM APIs—need to be driven by infrastructure as code so they can scale with any change to data. Functionalities that do not meet this fit should be minimized and deprioritized.

Policy-driven automation ensures staff know their roles and responsibilities for handling sensitive PII data. Identity and access management must mitigate risk by supporting policy enforcement at the point of access. Ideally, that happens before engineers, data scientists, and others use the data: that supports their know-your-customer and other guidelines and reduces the number of exceptions needed. It is critical to know which roles the different classification levels support—and do not support—before users submit requests for exception handling.

Using Unity Catalog to implement the above paradigm automates a significant part of meeting PII needs. Databricks PII Scanner enables automated metadata tagging of PII at the column level. By tracking those tags—combined with Teradata database-level PII classifications—enforcement points are established for automating PII policy exceptions, PII exceptions for Databricks Discover Permission Access, and alerting stakeholders for Databricks Export Task Access. What was perceived as, at best, a manual and resource-intensive process is expressed as an automated code solution.

3.1. Data Discovery and Metadata Management

Implementing Unity Catalog, a hat-landing service in Databricks, addresses the metadata management and PII identification requirements. It uses metadata tags and labels assigned to data assets and supports automatic enforcement of PII policies. Data Assets within Unity Catalog contain business and technical metadata, data lineage information, and security-related information. The data lineage graphs of Unity Catalog provide a PII data-flows visualisation. The metadata layer is populated through the integration with data sources, databases, and AI models within the Databricks environment. Integrated data-science and machine-learning models detect PII elements in semi-structured and structured data formats. A set of ingested and classified documents are used as training data to build a Machine Learning (ML) classifier model that automatically classifies the data set to predict the presence of PII information. The



metadata of new incoming data is processed and stored in the Unity Catalog data-management layer. Such metadata, together with the ML classification model output for each new ingested data set, provides a foundation for PII policy design and automation of subsequent workflows.

The business, technical, and security metadata of the Unity Catalog support the design and automation of enforcement steps related to PII policy definition and implementation. Tagging or labelling data assets provides visibility of PII assets and assists in preventing unnecessary data replicas and synonym data repositories. Exposure of real-world PII to untrusted or unauthorised users is the main concern of PII data regulation frameworks. By integrating Microsoft Active Directory Identity and Access Management with Databricks, fine-grained access to PII data at and below different data-storing locations can be managed. The managed Microsoft Active Directory groups define Databricks workspace roles for Databricks work area Approval processes that consider end-user data-consumption objectives, PII data presence in the dataset, data governance roles, and regulatory requirements can be automated to eliminate governance bottlenecks and speed-up business operations.

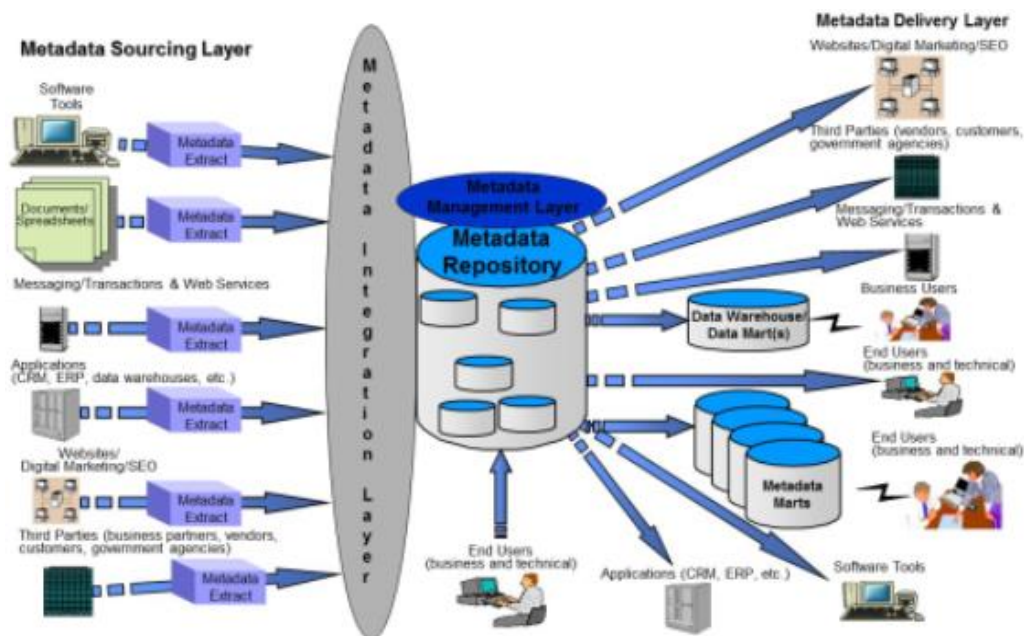


Fig 3: Metadata Management

3.2. Access Control and Policy Language

Data governance and compliance policies are often complex and require numerous supporting activities and processes. Organizations develop these policies to meet legal and regulatory requirements, protect customer trust, prevent financial loss, reduce risk, and adhere to market best practices. These policies are difficult to realize in practice because of the large number of Key Performance Indicators (KPIs) and the associated activities. For operational flexibility, organizations assign responsibility for these KPIs across organizational units rather than a single data governance function or unit. Therefore, information lifecycle management and PII compliance are highly distributed across the organization, leading to increased operational complexity.

Automation of information lifecycle management and PII compliance can greatly reduce the manual effort required. Automated policy checks also enable integration of multiple policies across service lines. The justification for PII identifiability for cloud services can be scanned and automatically included in the work products, and automated report generation reduces the need for manual compilation of reports. Policy execution failure can automatically trigger exception management to seek approval for data usage. Data owners need to continuously maintain a business view of the PII relevant to the data under their ownership, and the identification of de-ascribed data should be incorporated in the production data catalog.

Equation 3: k-anonymity equation
Step 1: Define an equivalence class



Let E_j be the set of records sharing the same quasi-identifier values.

Example:

all records with the same generalized age, ZIP prefix, and gender.

Step 2: Count class size

Let $|E_j|$ be the number of records in that class.

Step 3: State anonymity requirement

A dataset satisfies k -anonymity iff every equivalence class has at least k records:

$$|E_j| \geq k \forall j$$

Step 4: Why this matches the text

If one person is in a group of size k , that person is indistinguishable from $k - 1$ others.

Final derived equation

$$\boxed{\forall j, |E_j| \geq k}$$

IV. OBJECTIVE OF THE STUDY

Effective data governance involves roles, responsibilities, accountability, and control over data assets and resources, enabling compliance with legal and regulatory requirements. Integrating a data discovery tool with metadata management and automation workflows enables organizations to automatically classify sensitive personally identifiable information (PII), manage access control based on predefined policies, execute data masking or de-identification techniques, and log the data access requests and decisions taken, including approval and denial workflows.

Achieving such comprehensive end-to-end automation requires integrating an identity and access management (IAM) management platform to assign different roles either based on existing IAM integration or an analytical model. Any sensitive data identified through data discovery can automatically trigger authentication. The Vault key management system then seamlessly automates the recording of encryption keys used for masking so that the organization does not unnecessarily derive or store the de-identified values of PII.

4.1. Policy-Driven Automation Workflows

Multiple overlapping data governance frameworks govern data within an enterprise. These frameworks require automation mechanisms to reduce complexity, resource consumption, and errors. Such automations manage aspects including data discovery and inventory, metadata exposure and management, and operation exposure to protect sensitive information. A data governance system such as Databricks Unity Catalog supports a unified environment for data and AI capabilities across clouds with fine-grained access control. Policy-driven automations operate on the exposed metadata and Enterprise Data Catalog (edc) of the system. Techniques for automated operations include Data and Pseudonymization eNabling De-identification and Anonymization (DND). PII data identification and cataloging guarantee that PII-associated metadata and information are available throughout the organization. The system identifies PII status, makes PII policy and process assignments based on the PII exposure level, and supports a PII exception process.

Data governance frameworks and structures allocate governance roles across the organization. These roles address compliance with legal requirements, governing data regulation, and data protection policies. Role-based segregation of duties supports appropriate checks and balances in critical decision-making processes. A data governance system such as Databricks Unity Catalog secures PII data by identifying the data owner, providing approvals for elevated access roles, and enforcing separation of duties through structured approval processes. The engine explicitly enables data de-identification and anonymization capabilities at PII data exposure levels.

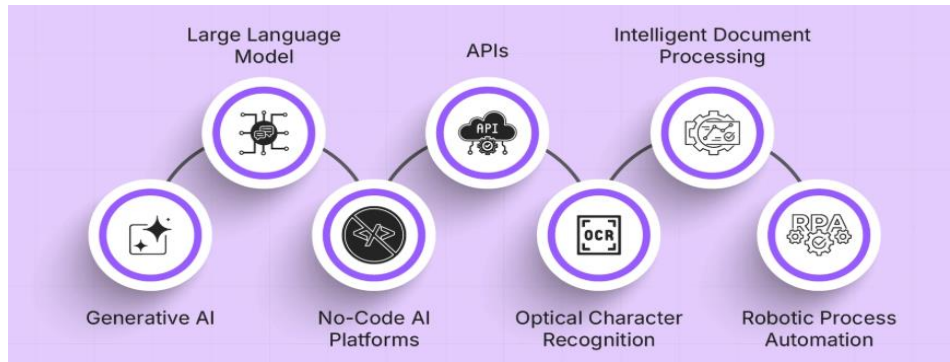


Fig 4: Technologies Used in AI Workflow Automation

4.2. Identity and Access Management Integration

Data governance and Personally Identifiable Information (PII) compliance processes span the entire ecosystem and overflow from data producers to consumers. Leveraging Identity and Access Management (IAM) solutions can help organizations address these challenges more effectively by facilitating discovery, accountability, prevention of unwanted information leaks, and speeding up PII-related processes such as remediation, management of legal holds, and audit logging.

When sensitive information has been identified and appropriate remediating actions completed, IAM integrations with data environments such as Unity Catalog can help enforce the desired policies through hidden columns or automatic data masking, anonymization, or de-identification. When teams process sensitive information for sensitive purposes within data environments, these integrations can help make compliance delegation easier and machine learning/artificial intelligence (ML/AI) models shallower and safer by easily identifying the presence of sensitive attributes during training.

Equation 4: l-diversity equation

Step 1: Let $S(E_j)$ be the set of distinct sensitive values in class E_j

Example:

diagnosis values inside one anonymized group.

Step 2: Count distinct values

$$|S(E_j)|$$

Step 3: Require at least l different sensitive values

$$|S(E_j)| \geq l$$

Final derived equation

$$\forall j, |S(E_j)| \geq l$$

Entropy form

A stronger form is entropy l-diversity:

$$-\sum_{s \in S(E_j)} p_j(s) \log p_j(s) \geq \log l$$

where $p_j(s)$ is the fraction of records in class E_j having sensitive value s .

Step-by-step for entropy form

1. For each sensitive value s , compute probability

$$p_j(s) = \frac{\text{count of } s \text{ in } E_j}{|E_j|}$$

2. Compute entropy



$$H(E_j) = - \sum_s p_j(s) \log p_j(s)$$

3. Impose diversity threshold

$$H(E_j) \geq \log l$$

The simpler derived equation usually expected is:

$$|S(E_j)| \geq l$$

V. METHODOLOGY

The identification and cataloging of personally identifiable information (PII) constitute the first step toward compliance with data-privacy legislation. A taxonomy of PII is created, focusing on its specific attributes of identification, linkage, and de-identification. PII is classified into three categories based on re-identification risk, identity exposure, type of consent for data usage, and usage of additional information to reduce the risk of re-identification. De-identification and anonymization techniques are reviewed with examples specific to each type of PII. These foundation concepts are essential for the successful implementation of policy-driven automation in data ecosystems.

With the rapid adoption of cloud technology, organizations are migrating their hybrid data estate into the cloud. While this enables organizations to benefit from a wide range of services, it also poses new challenges for data governance and privacy compliance. Organizations must ensure that sensitive data is being scanned, classified, and managed according to well-established, inclusive, and automated processes. Inadequate data governance can expose organizations to significant financial risk and reputational loss. A policy-driven approach to automation helps in instilling the desired roles, responsibilities, and accountability for sensitive data across its lifecycle.

5.1. PII Identification and Cataloging

Enrich the existing metadata layer of the data estate with additional information necessary to identify and manage PII reliably. Discover additional sources of sensitive data and perform one-time filling of the corresponding metadata fields. Change the metadata of existing data assets to reflect new information. Apply appropriate categories to tables, views, files, columns, and cells to reflect PII content, sensitivity level, de-identification status, and access restrictions.

The data governance function must be empowered to steward the metadata layer of the new platform. Business experts must identify the most critical sources of PII that need to be managed manually. The metadata data model should require business owners or security officers to mark tables, files, columns, views, and cells containing PII and indicate their sensitivity level. It should also provide means for indicating the status of de-identification associated with data that can be further processed. Finally, in a business operating model offering PII data to machine learning teams, the management of unstructured data sources should be extended beyond detection and warnings to manual PII cataloging.

Equation 5: t-closeness equation

Step 1: Define global sensitive distribution

Let $P(S)$ be the overall dataset distribution of the sensitive attribute.

Step 2: Define class-level distribution

Let $P(S | E_j)$ be the sensitive-value distribution inside equivalence class E_j .

Step 3: Measure distance between the two distributions

Use a distribution distance $d(\cdot, \cdot)$, commonly Earth Mover's Distance.

$$d(P(S | E_j), P(S))$$

Step 4: Impose closeness threshold

For t-closeness,

$$d(P(S | E_j), P(S)) \leq t$$

5.2. De-identification and Anonymization Techniques

The de-identification and anonymization techniques described in the pioneers' analysis are relevant to any organization processing PII, such as Healthcare and Life Sciences, Banking and Financial Services, Retail, Consumer Products, and Manufacturing, including AI-driven environments. PII consists of data or information that can be used to identify a particular individual, and strict compliance regulations regulate its usage. Several approaches have been adopted to



protect such sensitive data in data analytics and machine learning while satisfying GDPR and CCPA privacy regulations, which enforce security by default. A formal risk-based approach to the usage of these data has been established previously. Their risks are divided into different categories, and data usage is processed accordingly, rejecting some types of usage a priori if it is considered too critical. The GDPR-CCPA regulatory impact can also be monitored using these risk signals.

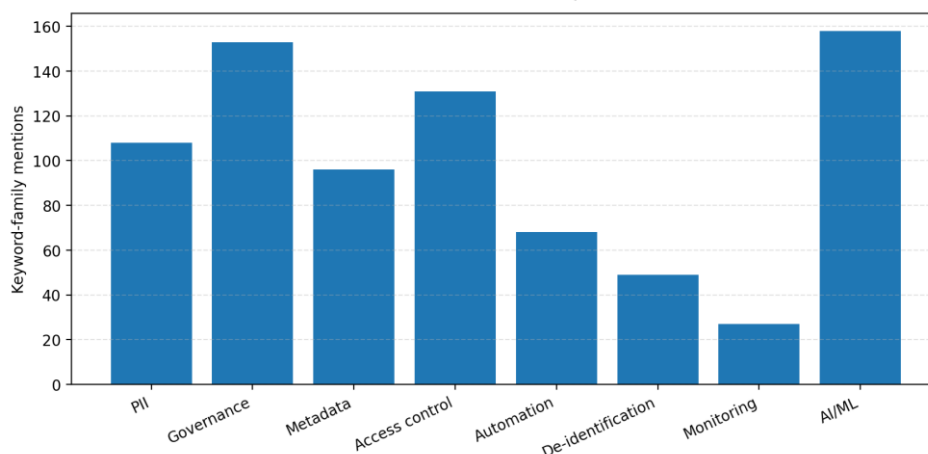
The most common de-identification and anonymization techniques that can be used in the DataBricks environment are outlined as follows. **Generalization** replaces a precise value with an imprecise value and is used to avoid coding with a k-anonymity property. For example, original value: age 27 becomes generalized age [20–30]. **K-anonymity** is achieved by making changes like generalization, suppression, and adding noise to a PII subset or group of records. Any individual in that group is indistinguishable from at least k-1 others in relation to some sensitive attributes. Other common techniques used in combination are **l-diversity** and **t-closeness**, which consider sensitive attribute values. For example, l-diversity adds clutter for sensitive categories, and t-closeness considers distribution. **Data perturbation** adds noise to continuous value attributes, retaining overall value of data. **Data-swapping**, **data-sampling**, and **masking** are other techniques. Data that does not need to be reversible can use simple masking techniques. For very sensitive data, like medical records, anonymization using data-relational-grid-aware synthetic data generation is highly recommended.

VI. RESULT

Role-based access control is a method for regulating access based on users' roles within a system. A role is defined according to authority and responsibility within an organization. Moreover, roles can extend beyond a particular organization and can be assigned to applications and services for controlling security operations like Single Sign-On. Access and error logs track events associated with virtually all information-processing systems. This log data provides important information regarding who accessed private data and whether compliance policies are being followed or, at least warranted. Such logs can detect erroneous and unintended access to private data and therefore can support accountability mechanisms.

Policy compliance requires that new policies designed to enforce compliance and relevant exceptions are continuously monitored to verify that they are correctly enforced and that the roles performing the actions specified in the exceptions are complying with the intended compliance policy together with the exceptions. Violation of this monitoring can then be further investigated. Monitoring is also required in crating ad-hoc de-identification or anonymization processes of private data when the data is used outside of the defined access control mechanisms.

Derived Bar Chart 1 - Theme Emphasis in the Article



6.1. Roles, Responsibilities, and Accountability

Roles, responsibilities, and accountability for data governance across Identity and Access Management (IAM) systems and data handling systems are not new subjects in a corporate environment; they are precious for the accomplished execution of compliance and security regulations. However, an organization does not give value to them until it reaches the stage where flexibility of access to resources, especially personal data, meets the boundary limits of compliance and security.



Policy-based identification of attributes in IAM, debugging of sensitive information from data handling systems, for example, DataBricks Unity Catalog facilitating the ability to mark a column as sensitive, prevent the necessary civil work in preparation for an audit. Organizations can proactively manage who is processing sensitive data and for what reasons by marking that data as sensitive, enabling the logging of activity and triggering workflows to notify impacted or responsible organizations or teams when new sensitive data is introduced, modified or deleted.

6.2. Policy Enforcement and Exceptions Management

The success of any compliance initiative, including PII, depends on the definition of appropriate policies, including the necessary roles and responsibilities associated with data understanding, controls, and audit. A distinction is required between data stewards and data custodians, with the former being primarily responsible for understanding and maintaining the data and the latter, often IT, responsible for technical system administration. A council of trust can manage compliance and exceptions beyond local reporting lines. Roles supported by reporting tools provide a check on segregated duties.

The policies governing data remaining enable non-PII datasets and/or tables used for data-based decisions despite containing sensitive information to exist when adequately protected by masking at runtime and approved for production use. All exceptions or releases to PII data, generally outside of local interest, must be approved by senior management, with a dedicated PII process foreseen and working to formalize.

VII. CASE STUDY: IMPLEMENTING UNITY CATALOG FOR PII COMPLIANCE

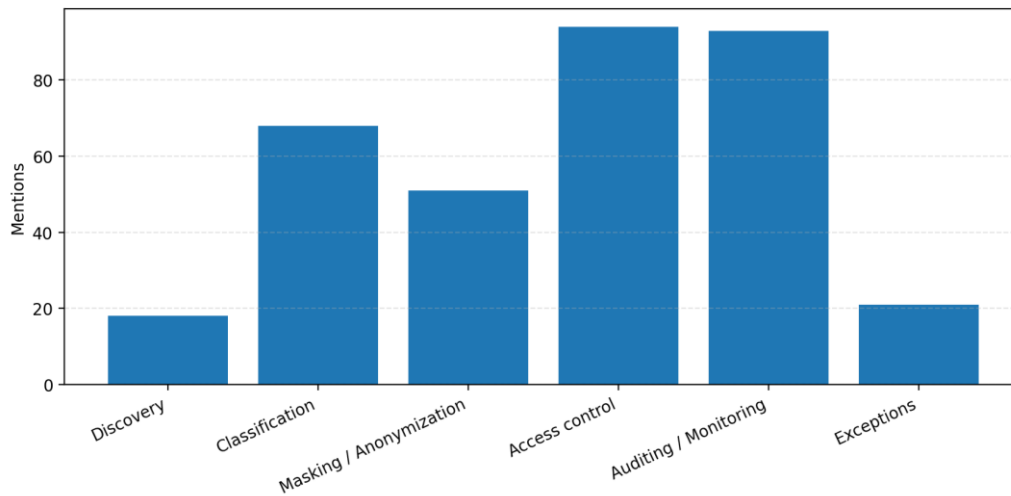
The implementation of Unity Catalog in a data ecosystem is examined by focusing on the PII policy objective, which automates the discovery, classification, isolation, and management of PII data throughout its life cycle.

A modern data ecosystem typically consists of integrated on-premise and cloud systems, supporting both batch and real-time data streams. The PII policy-supporting environment is represented diagrammatically, with its major components highlighted. Logs from various sources, such as access management, CloudTrail, or Apache Atlas, are gathered into a unified monitoring solution like DataDog for identity-aware access profiling. PII policies are defined in a suitable language having an expressive semantics and mapping functions for generating refined access control clauses across the data ecosystem. Data flows are monitored using tools such as Applinks, and the functioning of the policy-supporting environment is demonstrated by applying it to detect and mitigate data-sampling bias.

The formulation of rules for proactive management of data access, usage, and movement is a complex task, demanding large overhead, improved turnaround time, and inherent risks if done manually. The deployment of Unity Catalog enables a unified view of data assets, centralized policy management and enforcement, controlled data sharing, audit capabilities, tracking of sensitive information, and operational transparency. PII information from relevant databases has been collected into the catalog along with details of existing primary data producers and consumers, serving as the basis for role-based services matching and PII compliance on-ground support instructions.



Derived Bar Chart 2 - Compliance Control Points Mentioned in the Article



7.1. System Landscape and Data Flows

Data governance objectives can be pursued by deploying the Unity Catalog in the Databricks environment, which is ideally suited to data engineering, data science, and ML workloads. Unity Catalog provides end-to-end visibility through centralized metadata management, thereby addressing both data-discovery and governance requirements. PII policies can be automated via a combination of predefined, custom, and domain-specific classification rules.

The implementation example presented uses a heterogeneous system landscape comprising a Databricks environment on Microsoft Azure, a Microsoft Power BI dashboard, an Azure Data Lake Store, and an Azure SQL relational database. The AI model in the Databricks environment predicts whether a charge-off will occur on financial products. Data from the SQL database is joined with sensitive test and training data stored in the Azure Data Lake Store. A test candidate with the variable Target is predicated. The Power BI dashboard fetches the result of the prediction model for business users as part of a data science operation. The business concern is to protect sensitive customer data while creating an effective prediction model. The Unity Catalog verifies whether sensitive data is accessed in line with PII policy.

7.2. Policy Design and Automation Steps

Three standard organizational policies were designed to govern PII and the associated data-sharing de-identification approach within the Databricks ecosystem. These policies enabled end-to-end automation of policy enforcement, including exception handling for data-sharing rules involving de-identified data. Policies were defined as follows and the steps were presented in Figure 19.

1. PII Identification and Cataloguing: PII in the data ecosystem was identified and catalogued using DTL, along with respective business rules for de-identification and data-sharing approvals. Unity Catalog was used for PII asset management, DTL Metadata was stored in Delta Lake, and Databricks Workflows orchestrated the solution.

2. De-identification and Anonymization Techniques: Data stewards for PII were determined and de-identification/anonymization approaches were defined. Sensitive attributes were labelled in Unity Catalog and policy implementations were defined for data sharing. The Identity and Access Management (IAM) registration of stewards was kept up-to-date.

3. PII Policy Exceptions Management: A governance policy was established to approve exceptions requested by data consumers. Data-sharing requests that required the use of a specific sensitive attribute were handled through the PII policy exception process. Standard templates were created for exception requests and approvals to ease fulfilment.

VII. CONCLUSION

The ultimate goal of modern businesses is to give customers the finest experience possible in the least amount of time. To do this efficiently and effectively, many companies are now adopting new technologies, including deep learning-based chatbots and virtual assistants. But what about ensuring the protection of their customers' privacy when these



products use organizational data? Companies that have privacy-sensitive data collected in production systems and create data flows that lead the data into a sandbox or prototype environment in order to build a model should be careful with those automations to protect the privacy of their customers. Failure to protect private data, especially Personally Identifiable Information (PII), has serious consequences, from reputational damage to legal implications. The Project is aimed at providing PII protection to the design and implementation of automated data flows that support the use of AI and deep learning algorithms while determining the usage of such flows.

Candidates of systems that support new technologies should implement a centralized catalog that is capable of monitoring the data based on well-defined privacy rules. The solution is to automate PII identification and cataloging, execute de-identification and anonymization when needed, inform involved data users about PII presence on accessible datasets and provide the possibility of exception management workflow. Unity Catalog from Databricks is proposed to be the catalog system for implementation. Candidate design includes the organization of system landscape and data flows in order to use an IAM solution to determine users and their roles in the organization, description of privacy rules into Unity Catalog, formalization of PII identification and cataloging processes, and description of the automation flow to be created, comprising PII detection, usage of corresponding data policy, communication to data users and the exception management workflow.

Capability	How the article frames it	Primary mechanisms mentioned	Operational outcome
PII discovery & cataloging	Foundational step for compliance in AI-driven ecosystems	Unity Catalog metadata, tags, scanners, business rules	Visibility of sensitive assets and reduced blind spots
Access control	Enforced at point of access using identity-aware roles	IAM integration, RBAC, fine-grained permissions	Controlled use of sensitive datasets
De-identification	Applied when analytics needs must coexist with privacy requirements	Masking, anonymization, generalization, k-anonymity	Lower re-identification risk
Workflow automation	Manual approvals and checks are converted into repeatable flows	Policy-driven workflows, Databricks Workflows, notifications	Lower effort, faster turnaround, fewer bottlenecks
Audit & monitoring	Needed to verify policy execution and investigate violations	Logs, lineage, exception tracking, alerting	Accountability and evidence for compliance

Table : Capability matrix derived from the paper's main themes.

IX. LIST OF IMPORTANT REFERENCES

These references provide important background and support for the more detailed sections of the research. Machine learning has many challenges, one of which is the preparation of data for the model's training phase. The principle for preparing this data is that it should represent the production scenarios where the predictive model will be implemented. For this data preparation, one group of techniques, referred to as data transformation techniques, is used to modify the original data. Here, data transformation techniques together with de-identification and anonymization techniques are examined for privacy and security of personal data finding their importance in machine learning models. In this paper specific focus has been given on implementing data transformation techniques in model training phase.

Machine learning has many challenges, one of which is the preparation of data for the model's training phase. The principle for preparing this data is that it should represent the production scenarios where the predictive model will be implemented. For this data preparation, one group of techniques, referred to as data transformation techniques, is used to modify the original data. Here, data transformation techniques together with de-identification and anonymization



techniques are examined for privacy and security of personal data finding their importance in machine learning models. In this paper specific focus has been given on implementing data transformation techniques in model training phase.

REFERENCES

- [1] Aitha, A. R. (2023). CloudBased Microservices Architecture for Seamless Insurance Policy Administration. *International Journal of Finance (IJFIN)-ABDC Journal Quality List*, 36(6), 607-632.
- [2] Dwaraka Nath Kummari, Srinivasa Rao Challa, "Big Data and Machine Learning in Fraud Detection for Public Sector Financial Systems," *International Journal of Advanced Research in Computer and Communication Engineering (IJARCCE)*, DOI: 10.17148/IJARCCE.2020.91221
- [3] Sheelam, G. K., & Nandan, B. P. (2022). Integrating AI And Data Engineering For Intelligent Semiconductor Chip Design And Optimization. *Migration Letters*, 19, 2178-2207.
- [4] Mangalampalli, B. M. (2023). AI-Driven Anomaly Detection in Healthcare Claims Data: A Business Intelligence Perspective. *Journal of Rare Cardiovascular Diseases*.
- [5] Mukesh, A., & Aitha, A. R. (2021). Insurance Risk Assessment Using Predictive Modeling Techniques. *International Journal of Emerging Research in Engineering and Technology*, 2(4), 68-79.
- [6] Palanichamy, R. S. T. (2023). AI and data governance: Enhancing security, privacy, and accountability. *International Journal on Science and Technology*, 14(1), 1–10
- [7] Kolla, S. K. (2023). Explainable AI and ML Models for Transparent Clinical Decision Support. *Journal for ReAttach Therapy and Developmental Diversities*, 6, 2444-2460.
- [8] Meda, R. End-to-End Data Engineering for Demand Forecasting in Retail Manufacturing Ecosystems.
- [9] Gadi, A. L. , Gadi, A. L. Kannan, S. , Kannan, S. Nandan, B. P. , Nandan, B. P. Komaragiri, V. B. , & Komaragiri, V. B. (2021). Advanced Computational Technologies in Vehicle Production, Digital Connectivity, and Sustainable Transportation: Innovations in Intelligent Systems, Eco-Friendly Manufacturing, and Financial Optimization. *Universal Journal of Finance and Economics*, 1(1), 87-100. <https://doi.org/10.31586/ujfe.2021.1296>.
- [10] Inala, R. Advancing Group Insurance Solutions Through Ai-Enhanced Technology Architectures And Big Data Insights.
- [11] Kannan, S., Nuka, S. T., Pamisetty, V., Gadi, A. L., Krishna, H., & Koppolu, R. ENHANCING AGRICULTURAL EQUIPMENT AND MEDICAL DEVICES Pamisetty, V. (2020). Optimizing tax compliance and fraud prevention through intelligent systems: The role of technology in public finance innovation. Available at SSRN 5250796.
- [12] Kummari, D. N., & Burugulla, J. K. R. (2023). Decision Support Systems for Government Auditing: The Role of AI in Ensuring Transparency and Compliance. *International Journal of Finance (IJFIN)-ABDC Journal Quality List*, 36(6), 493-532.
- [13] Kalisetty, S., & Singireddy, J. (2023). Optimizing Tax Preparation and Filing Services: A Comparative Study of Traditional Methods and AI Augmented Tax Compliance Frameworks. Available at SSRN 5206185.
- [14] Adusupalli, B., Singireddy, S., & Pandiri, L. Implementing Scalable Identity and Access Management Frameworks in Digital Insurance Platforms. *International Journal of Advanced Research in Computer and Communication Engineering (IJARCCE)*, DOI, 10.
- [15] Segireddy, A. R. (2022). Terraform and Ansible in Building Resilient Cloud-Native Payment Architectures. *International Journal of Intelligent Systems and Applications in Engineering*, 10, 444-455.
- [16] Gottimukkala, V. R. R. (2020). Energy-Efficient Design Patterns for Large-Scale Banking Applications Deployed on AWS Cloud. *power*, 9(12).
- [17] Garapati, R. S., & Kanna, S. R. A Digital Twin-Enabled Predictive Maintenance Framework Leveraging Multi-Agent Reinforcement Learning and Industrial IoT Data.
- [18] Pamisetty, V., Dodda, A., Lakarasu, P., Singireddy, J., & Challa, K. (2022). Optimizing Digital Finance and Regulatory Systems Through Intelligent Automation, Secure Data Architectures, and Advanced Analytical Technologies. *Secure Data Architectures, and Advanced Analytical Technologies* (December 10, 2022).
- [19] Nasiri, S., et al. (2023). A systematic review of big data stream processing frameworks and applications. *Journal of Big Data*, 10(1), 67.
- [20] Mangalampalli, B. M. Intelligent Data Profiling for Healthcare Data Lakes Using AI-Enhanced Analytics.
- [21] Kolla, S. H. (2023). Deep Learning–Driven Retrieval-Augmented Generation for Enterprise ITSM Automation: A Governance-Aligned Large Language Model Architecture. *Journal of Computational Analysis and Applications*, 31(4).
- [22] Singireddy, J. (2022). Leveraging Artificial Intelligence and Machine Learning for Enhancing Automated Financial Advisory Systems: A Study on AIDriven Personalized Financial Planning and Credit Monitoring. *Mathematical Statistician and Engineering Applications*, 71(4), 16711-16728.



- [23] Amistapuram, K. Energy-Efficient System Design for High-Volume Insurance Applications in Cloud-Native Environments. *International Journal of Innovative Research in Electrical, Electronics, Instrumentation and Control Engineering (IJIREEICE)*, DOI, 10.
- [24] Mahesh Recharla, (2020), "Targeted Gene Therapy for Spinal Muscular Atrophy: Advances in Delivery Mechanisms and Clinical Outcomes", *International Journal of Science and Research (IJSR)*, 9(12), 1921-1934. <https://dx.doi.org/10.21275/SR20126161624>, <https://www.ijsr.net/getabstract.php?paperid=SR20126161624>
- [25] Kulkarni, A. R., Kumar, N., & Rao, K. R. (2023). Big data analytics and monitoring frameworks for scalable data pipelines. *Big Data Mining and Analytics*, 6(2), 139–153.
- [26] Botlagunta Preethish Nandan, "Data Analytics-Driven Approaches to Yield Prediction in Semiconductor Manufacturing," *International Journal of Innovative Research in Electrical, Electronics, Instrumentation and Control Engineering (IJIREEICE)*, DOI 10.17148/IJIREEICE.2021.91217.
- [27] Garapati, R. S. (2023). Optimizing Energy Consumption in Smart Build-ings Through Web-Integrated AI and Cloud-Driven Control Systems.
- [28] Chowdhury, R. H. (2021). Cloud-based data engineering for scalable business analytics solutions: designing scalable cloud architectures to enhance the efficiency of big data analytics in enterprise settings. *Journal of Technological Science & Engineering (JTSE)*, 2(1), 21-33.
- [29] Vamsee Pamisetty, Lahari Pandiri, Sneha Singireddy, Venkata Narasareddy Annapareddy, Harish Kumar Sriram. (2022). Leveraging AI, Machine Learning, And Big Data For Enhancing Tax Compliance, Fraud Detection, And Predictive Analytics In Government Financial.
- [30] Gottimukkala, V. R. R. (2021). Digital Signal Processing Challenges in Financial Messaging Systems: Case Studies in High-Volume SWIFT Flows.
- [31] Aitha, A. R. (2023). Cloud-Native Big Data AI/ML Framework for Risk Intelligence and Fraud Control in Banking and Insurance Ecosystems. Available at SSRN 6157967.
- [32] Sheelam, G. K., & Nandan, B. P. (2021). Machine Learning Integration in Semiconductor Research and Manufacturing Pipelines. *International Journal of Advanced Research in Computer and Communication Engineering (IJARCC)*, DOI, 10.
- [33] Chakilam, C., Suura, S. R., Koppolu, H. K. R., & Recharla, M. (2022). From Data to Cure: Leveraging Artificial Intelligence and Big Data Analytics in Accelerating Disease Research and Treatment Development. *Journal of Survey in Fisheries Sciences*. <https://doi.org/10.53555/sfs.v9i3.3619>.
- [34] Nagabhyru, K. C. (2023). Accelerating Digital Transformation with AI Driven Data Engineering: Industry Case Studies from Cloud and IoT Domains. *Educational Administration: Theory and Practice*, 29(4), 5898-5910
- [35] Bonawitz, K., et al. (2023). Secure aggregation for federated learning. Google Research.
- [36] Yandamuri, U. S. (2022). Big Data Pipelines for Cross-Domain Decision Support: A Cloud-Centric Approach. *International Journal of Scientific Research and Modern Technology (IJSRMT)*.
- [37] Davuluri, P. N. Integrating Artificial Intelligence into Event-Driven Financial Crime Compliance Platforms.
- [38] Gottimukkala, V. R. R. (2023). Privacy-Preserving Machine Learning Models for Transaction Monitoring in Global Banking Networks. *International Journal of Finance (IJFIN)-ABDC Journal Quality List*, 36(6), 633-652.
- [39] Dwaraka Nath Kummari,. (2022). Machine Learning Approaches to Real-Time Quality Control in Automotive Assembly Lines. *Mathematical Statistician and Engineering Applications*, 71(4), 16801–16820. Retrieved from <https://philstat.org/index.php/MSEA/article/view/2972>
- [40] Goutham Kumar Sheelam. (2022). Reconfigurable Semiconductor Architectures For AI-Enhanced Wireless Communication Networks. *Kurdish Studies*, 10(2), 1027–1040. <https://doi.org/10.53555/ks.v10i2.3867>.
- [41] Pamisetty, A. (2022). Big Data can Generate Major Opportunities for Manufacturing Supply Chains. *International Journal of Scientific Research and Modern Technology*, 1(12), 238–251. <https://doi.org/10.38124/ijsrmt.v1i12.1186>
- [42] Yandamuri, U. S. (2021). A Comparative Study of Traditional Reporting Systems versus Real-Time Analytics Dashboards in Enterprise Operations. *Universal Journal of Business and Management*
- [43] Garapati, R. S. (2022). AI-Augmented Virtual Health Assistant: A Web-Based Solution for Personalized Medication Management and Patient Engagement. Available at SSRN 5639650.
- [44] Inala, R. Designing Scalable Technology Architectures for Customer Data in Group Insurance and Investment Platforms.
- [45] Kolla, S. H. (2021). Rule-Based Automation for IT Service Management Workflows. *Online Journal of Engineering Sciences*, 1(1), 1-14.
- [46] Segireddy, A. R. (2020). Cloud Migration Strategies for High-Volume Financial Messaging Systems.
- [47] Yandamuri, U. S. (2023). An Intelligent Analytics Framework Combining Big Data and Machine Learning for Business Forecasting. *International Journal Of Finance*, 36(6), 682-706.
- [48] Singireddy, J. (2023). Finance 4.0: Predictive analytics for financial risk management using AI. *European Journal of Analytics and Artificial Intelligence (EJAAI)* p-ISSN, 3050-9556.



- [49] Somasundaram, P. (2023). Improving real-time job monitoring for cloud-based data pipelines. *International Journal of Computer Engineering and Technology*, 14(3), 39–47.
- [50] Davuluri, P. N. (2020). Event-Driven Architectures for Real-Time Regulatory Monitoring in Global Banking.
- [51] Kolla, T. (2023). Predictive ETL Failure Detection in Healthcare Data Pipelines Using Anomaly Detection Algorithms. *International Journal of Medical Toxicology & Legal Medicine*.
- [52] Nagabhyru, K. C. (2023). From Data Silos to Knowledge Graphs: Architecting CrossEnterprise AI Solutions for Scalability and Trust. Available at SSRN 5697663.
- [53] Recharla, M., & Chitta, S. AI-Enhanced Neuroimaging and Deep Learning-Based Early Diagnosis of Multiple Sclerosis and Alzheimer's.
- [54] Aiswarya, K., Reddy, P., & Kumar, V. (2023). Fault detection and mitigation strategies in data pipeline systems. *International Journal of Data Engineering*, 14(1), 22–34.
- [55] Botlagunta, P. N., & Sheelam, G. K. (2020). Data-Driven Design and Validation Techniques in Advanced Chip Engineering. *Global Research Development (GRD) ISSN*, 2455-5703.
- [56] Meda, R. (2020). Designing Self-Learning Agentic Systems for Dynamic Retail Supply Networks. *Online Journal of Materials Science*, 1(1), 1-20.
- [57] Valiki, D., & Kummari, D. N. (2021). Rule-Based Decision Systems for the Automation of Audit Sampling. *International Journal of Emerging Trends in Computer Science and Information Technology*, 2(4), 105-114
- [58] Mangala, N. (2021). CI/CD Pipeline Automation for Enterprise Data Artifacts Using Azure DevOps. *Universal Journal of Business and Management*, 1(1), 1-18. <https://doi.org/10.31586/ujbm.2021.1363>
- [59] Nagubandi, A. R. (2023). Advanced Multi-Agent AI Systems for Autonomous Reconciliation Across Enterprise Multi-Counterparty Derivatives, Collateral, and Accounting Platforms. *International Journal of Finance (IJFIN)-ABDC Journal Quality List*, 36(6), 653-674
- [60] Amistapuram, K. (2022). Fraud Detection and Risk Modeling in Insurance: Early Adoption of Machine Learning in Claims Processing. Available at SSRN 5741982.
- [61] Mangala, N. (2022). Real-Time Data Quality Monitoring and Gating Frameworks in Cloud-Based Data Pipelines. *International Journal of Research and Applied Innovations*, 5(6), 8197-8219.
- [62] Nasiri, S., Rahmani, A. M., & Rezaei, M. (2023). A systematic review of big data stream processing frameworks and applications. *Journal of Big Data*, 10(1), 67.
- [63] Inala, R. (2021). A New Paradigm in Retirement Solution Platforms: Leveraging Data Governance to Build AI-Ready Data Products. *Journal of International Crisis and Risk Communication Research*, 286-310.
- [64] Pamisetty, A. (2021). A comparative study of cloud platforms for scalable infrastructure in food distribution supply chains.
- [65] Malempati, M., Pandiri, L., Paleti, S., & Singireddy, J. (2023). Transforming financial and insurance ecosystems through intelligent automation, secure digital infrastructure, and advanced risk management strategies. *Jeevani, Transforming Financial And Insurance Ecosystems Through Intelligent Automation, Secure Digital Infrastructure, And Advanced Risk Management Strategies (December 03, 2023)*.
- [66] Pamisetty, A. (2022). Integrating Big Data, AI, and Financial Modeling in Cloud-Based Insurance and Banking Ecosystems. *AI, and Financial Modeling in Cloud-Based Insurance and Banking Ecosystems (December 05, 2022)*.
- [67] Sriram, H. K., ADUSUPALLI, B., Singireddy, S., & Malempati, M. (2021). Revolutionizing Risk Assessment and Financial Ecosystems with Smart Automation, Secure Digital Solutions, and Advanced Analytical Frameworks. *Murali, Revolutionizing Risk Assessment and Financial Ecosystems with Smart Automation, Secure Digital Solutions, and Advanced Analytical Frameworks (December 27, 2021)*.