



Policy-as-Code for AI Data Platforms: Enforcing Privacy, Lineage, and Access Control End-to-End

Gunda Vamshi Krishna

Department of Computer Science, Anurag University, Hyderabad, India

* Corresponding Author: **Gunda Vamshi Krishna**

Article Info

ISSN (online): 3107-3972

Volume: 01

Issue: 06

November-December 2024

Received: 10-09-2024

Accepted: 08-10-2024

Published: 06-11-2024

Page No: 228-232

Abstract

AI data platforms increasingly blend batch and streaming pipelines, lakehouse storage, feature stores, and model serving into a single production surface area. This convergence increases governance complexity: privacy constraints must survive transformations, lineage must remain provable across heterogeneous tools, and access control must be consistent from raw ingestion through feature creation to model inference. This paper proposes an end-to-end Policy-as-Code (PaC) reference architecture for AI data platforms that unifies (i) privacy policy enforcement, (ii) lineage capture and validation, and (iii) authorization and access control. The approach treats governance as versioned, testable, and deployable code artifacts that are compiled into platform-specific enforcement points while preserving a single source of truth. We define a policy taxonomy aligned to data and ML lifecycles, specify enforcement placement across the data plane and control plane, and introduce verification hooks that prevent “policy drift” between CI pipelines and runtime. The proposed design supports regulated domains and high-assurance requirements by integrating privacy budgeting, provenance-aware controls, and modern authorization languages to achieve consistent governance at scale.

DOI: <https://doi.org/10.54660/GMPJ.2024.1.6.228-232>

Keywords: Policy-as-Code, AI governance, data lineage, provenance, privacy, access control, lakehouse, authorization

1. Introduction

Modern organizations increasingly run ML workloads on shared data platforms where sensitive datasets, derived features, and trained models coexist and are reused across teams, which amplifies the blast radius of governance failures. ^[1]

In high-stakes domains such as clinical decision support, the operational value of AI is inseparable from the need to enforce privacy and compliance constraints across the full data-to-model pipeline. ^[2]

Governance frameworks for trustworthy AI emphasize lifecycle-wide risk management, motivating controls that are enforceable, auditable, and repeatable rather than ad hoc. ^[3]

In parallel, operational quality problems (including data defects and pipeline regressions) can directly degrade downstream model behavior, reinforcing the need to treat platform policies as first-class engineering artifacts rather than external documentation. ^[4]

Research question: How can an AI data platform enforce privacy, lineage, and access control consistently across ingestion, transformation, feature engineering, training, and serving—without duplicating logic across tools?

Contributions

1. A PaC reference architecture that maps platform lifecycle stages to concrete enforcement points.
 2. A unified policy taxonomy covering privacy constraints, lineage requirements, and authorization decisions.
 3. A compilation and deployment model that preserves a single policy source while targeting heterogeneous engines.
 4. An assurance model for provable lineage and policy drift prevention.
-

2. Background and Motivation

2.1. Data governance in AI platforms

Systematic analyses of data governance highlight recurring organizational challenges: fragmented ownership, inconsistent standards, weak auditability, and inadequate measurement of data quality controls. ^[1]

AI platforms intensify these issues because derived artifacts (features, embeddings, labels, trained weights) are both high-value and high-risk, yet they frequently inherit metadata and constraints imperfectly.

2.2. Why Policy-as-Code now

PaC operationalizes governance by expressing rules in machine-evaluable formats stored in version control, enabling review, testing, and automated deployment. ^[5]

A practical challenge is that AI platforms typically contain multiple decision surfaces (SQL engines, orchestration layers, model registries, serving gateways), so PaC must support distributed enforcement rather than a single monolithic gate.

2.3. End-to-end compliance economics

Organizations adopt automation when it reduces toil and shortens feedback loops; governance automation similarly requires evidence that enforcement can be integrated into delivery pipelines without blocking velocity. ^[9]

This motivates designs that (i) shift enforcement left into CI, (ii) preserve fast runtime checks, and (iii) emit audit-ready evidence.

3. Problem Statement and Requirements

3.1. Threats and failure modes

AI data platforms face governance failures that recur in practice: policy drift (runtime behavior diverges from policy assumptions validated in CI), lineage gaps (transformations occur without standardized capture, breaking audit trails), privilege escalation (inconsistent authorization across layers enables unintended access paths), and privacy leakage (downstream artifacts expose sensitive information without budgeted safeguards).

Privacy and security failures are increasingly framed as part of broader online safety and trust requirements in the information era. ^[19]

3.2. Requirements

We formalize end-to-end PaC requirements:

R1. Single source of truth: policies are authored once and deployed consistently across enforcement points.

R2. Compositional enforcement: policies apply to datasets, derived artifacts, and models, including tag propagation rules.

R3. Provable lineage: lineage capture is standardized and independently verifiable.

R4. Consistent authorization: access decisions are consistent across APIs, compute engines, and orchestration layers.

R5. Privacy-aware computation: privacy constraints include transformation rules and optional privacy budgeting for ML workloads. ^[16]

R6. Auditability: each enforcement decision yields evidence suitable for compliance reviews.

R7. Portability: policies remain stable while platform components evolve.

4. Policy Taxonomy for AI Data Platforms

We define three primary policy families and one supporting family.

4.1. Authorization and access control policies

These policies decide “who can do what” over resources such as datasets, tables, features, model artifacts, and endpoints. Surveys of distributed-system access control emphasize the need to combine multiple paradigms (RBAC, ABAC, relationship-based models) to match modern multi-tenant environments. ^[14]

For data platforms, ABAC is particularly natural because policy decisions often depend on sensitivity tags, intended purpose, and environmental context. ^[7]

4.2. Privacy policies

Privacy policies constrain collection and retention, allowed transformations (masking, aggregation), downstream re-use boundaries, and optional differential privacy (DP) budgeting for repeated training tasks.

DP introduces a scarce “privacy resource” that must be scheduled and tracked to prevent cumulative leakage across pipelines. ^[16]

4.3. Lineage and provenance policies

Lineage policies require that critical operations emit lineage metadata and that derived artifacts retain links to sources and transformations.

Cloud-scale provenance studies highlight that provenance is essential for investigation and accountability, but also difficult due to heterogeneity and evidence complexity. ^[6]

Security-focused provenance research shows that provenance records can support detection, accountability, and non-repudiation when designed as tamper-evident chains. ^[10]

4.4. Assurance and evidence policies

These policies govern how enforcement evidence is produced and retained, including attestations, signed logs, and minimum audit fields.

5. Reference Architecture: End-to-End Policy-as-Code

5.1. Architectural overview

The PaC architecture is organized into four layers:

Policy authoring layer: Policies stored in version control, with review and CI tests.

Policy compilation and packaging layer: Static checks (linting, type validation) and compilation into deployable bundles.

Distribution and reconciliation layer: Bundles are distributed to enforcement points and reconciled continuously to avoid drift.

Enforcement layer: Multiple policy decision points (PDPs) and policy enforcement points (PEPs) across the AI platform. A widely used operational pattern is to integrate PaC into cloud governance workflows so policy checks run during provisioning, deployment, and runtime validation. ^[13]

5.2. Lineage standardization as a backbone

To prevent lineage gaps, lineage must be expressed in a shared schema that multiple tools can emit and consume.

OpenLineage provides a practical specification direction by representing lineage as interoperable events and extensible facets. ^[8]

Operational experience shows that even large-scale platforms must invest in performance and usability improvements to make lineage actionable for developers and data owners. ^[27]

5.3. Authorization language considerations

Authorization policies require expressiveness, analyzability, and performance under high request volume.

Cedar is a modern authorization language designed for readable policies with strong analysis properties and efficient evaluation. ^[28]

6. Enforcement Placement Across the AI Lifecycle

6.1. Ingestion and landing zone

Goal: ensure data enters the platform with correct classification, retention, and ownership metadata.

Enforcement points: schema and metadata validation in CI for ingestion jobs; runtime checks in ingestion services (rejecting untagged sensitive fields); and immutable logging of acceptance decisions.

6.2. Transformation and lakehouse compute

Goal: enforce transformation constraints (masking, purpose limitation) and ensure lineage emission.

Enforcement points: SQL/ETL job preflight policies; compute runtime interceptors that require lineage events for each run; and tag propagation checks for derived datasets.

6.3. Feature engineering and reuse

Goal: prevent uncontrolled feature sharing and ensure features retain lineage and sensitivity labels.

Enforcement points: feature registration policies (required metadata, allowed consumers); access control checks for feature retrieval; and lineage completeness checks before “production” designation.

6.4. Training pipelines

Goal: ensure training datasets meet privacy and authorization constraints and produce audit evidence.

Privacy-aware scheduling becomes relevant when teams repeatedly train on sensitive streams and must enforce cumulative privacy bounds. ^[16]

Threat modeling and governance are strengthened when platforms also recognize adversarial ML risk categories, including poisoning and evasion, as part of the enforcement context. ^[17]

6.5. Model registry and deployment

Goal: ensure model artifacts carry provenance links and are deployable only if policy checks pass.

Enforcement points: registry admission policies (must include lineage pointers and training data identifiers); deployment gating policies (must include approval evidence for sensitive contexts); and serving gateway authorization policies (per-tenant and per-purpose).

7. High-Assurance Extensions: Confidential and “Honest” Computing

7.1. Confidential computing for data-in-use

For particularly sensitive workloads, PaC can be strengthened by executing critical computations in trusted execution environments (TEEs) and binding evidence to attestation outputs.

Industry guidance frames confidential computing as protecting “data in use” via hardware-based isolation and

attestation mechanisms. ^[20]

Academic reviews caution that confidential computing definitions and comparisons can be ambiguous, motivating explicit threat models and careful claims about guarantees. ^[12]

A broader literature surveys the evolution and open challenges of confidential computing, supporting the view that it is becoming a foundational security primitive for regulated analytics. ^[22]

7.2. Demonstrable lineage and verifiable policy claims

An emerging perspective is that data governance needs demonstrable, externally verifiable lineage and provenance so policy can move from principle-based statements toward enforceable rules. ^[11]

This supports an “assurance mode” where (i) lineage events are signed, (ii) policy bundles are hashed and attested, and (iii) audit evidence is tamper-evident.

8. Validation and Evaluation Approach

8.1. Evaluation dimensions

We recommend evaluating a PaC deployment on coverage (fraction of platform surfaces with enforceable policies), consistency (absence of conflicts across enforcement points), latency overhead (added time per policy decision at runtime), audit completeness (whether evidence is sufficient for policy obligations), and operational toil (reduction in manual reviews and exception handling).

8.2. Compliance automation in hybrid environments

Compliance-as-Code research in hybrid cloud environments supports the feasibility of expressing compliance logic as automation artifacts integrated into delivery workflows. ^[15]

Continuous compliance automation also extends to security control verification processes, aligning with PaC’s requirement for repeatable checks and evidence production. ^[23]

9. Discussion

9.1. Trade-offs

Centralization vs. autonomy: a single source of truth reduces inconsistency but requires governance around policy ownership and change management.

Expressiveness vs. analyzability: more expressive policies can be harder to prove correct; typed, analyzable policy languages mitigate this tension. ^[28]

Lineage fidelity vs. performance: lineage capture must be efficient enough to be mandatory, or teams will bypass it.

9.2. Practical adoption path

A pragmatic rollout strategy: enforce mandatory metadata and tagging at ingestion; require lineage events for production pipelines; introduce consistent authorization checks at access edges; and expand privacy policies to training and serving, including budgeting when needed. ^[16]

9.3. Cybersecurity alignment

Because AI platforms are now core infrastructure, governance must align to broader cybersecurity objectives for public utilities and smart infrastructure contexts. ^[22]

Bridging traditional information security and modern cybersecurity considerations remains essential when platform policies must satisfy both operational and regulatory constraints. ^[24]

10. Conclusion

Policy-as-Code offers a practical, engineering-native approach to enforce privacy, lineage, and access control end-to-end in AI data platforms. This paper presented a reference architecture and policy taxonomy that treat governance as versioned, testable, deployable artifacts distributed across multiple enforcement points. Standardized lineage, analyzable authorization languages, and privacy-aware controls collectively reduce policy drift, improve auditability, and support safe reuse of data and models in regulated and high-stakes environments.

References

- Bernardo BMV, Mamede HS, Barroso JMP, dos Santos VMPD. Data governance & quality management—innovation and breakthroughs across different fields. *J Innov Knowl.* 2024;9(4):100598. doi:10.1016/j.jik.2024.100598.
- Kacheru G. Revolutionizing healthcare: the role of artificial intelligence in clinical practice. *J Comput Anal Appl.* 2023;31(4):1546–1544. Available from: <https://eudoxuspress.com/index.php/pub/article/view/3270>
- National Institute of Standards and Technology (NIST). Artificial Intelligence Risk Management Framework (AI RMF 1.0). NIST AI 100-1; 2023. doi:10.6028/NIST.AI.100-1. Available from: <https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf>
- Gunda SK. Comparative analysis of machine learning models for software defect prediction. In: *Proc Int Conf Power Energy Control Transmission Systems (ICPECTS)*; 2024; Chennai, India. p. 1–6. doi:10.1109/ICPECTS62210.2024.10780167.
- Sayfan G. Policy as code and the Open Policy Agent. Cisco Blogs. 2022 Mar 21. Available from: <https://blogs.cisco.com/developer/policyascode01>
- Abiodun OI, Alawida M, Omolara AE, Alabdulatif A. Data provenance for cloud forensic investigations: security, challenges, solutions and future perspectives. *J King Saud Univ Comput Inf Sci.* 2022;34(10 Pt B):10217–10245. doi:10.1016/j.jksuci.2022.10.018
- Tall AM, Zou CC. A framework for attribute-based access control in processing big data with multiple sensitivities. *Appl Sci.* 2023;13(2):1183. doi:10.3390/app13021183
- Parreira C. Introduction to OpenLineage. Adaltas. 2024 Feb 19. Available from: <https://www.adaltas.com/en/2024/02/19/openlineage-introduction/>
- Kacheru G, Bajjuru R, Arthan N. The ROI of software automation: measuring time and cost savings. *Int J Commun Netw Inf Secur.* 2023;15(4):774–785.
- Pan B, Stakhanova N, Ray S. Data provenance in security and privacy. *ACM Comput Surv.* 2023;55(14s):323. doi:10.1145/3593294
- Guitton F. Honest computing: achieving demonstrable data lineage and provenance for driving data and process-sensitive policies. *Data Policy Proc.* 2024 Dec 26. Available from: <https://www.cambridge.org/core/journals/data-and-policy/article/honest-computing-achieving-demonstrable-data-lineage-and-provenance-for-driving-data-and-process-sensitive-policies/B1F021006F26F202A262051002B50A26>
- Sardar MU, Fetzter C. Confidential computing and related technologies: a critical review. *Cybersecurity.* 2023;6:10. doi:10.1186/s42400-023-00144-1
- Timpone A, Banwart D. A practical guide to getting started with policy as code. AWS Infrastructure & Automation Blog. 2024 Dec 9. Available from: <https://aws.amazon.com/blogs/infrastructure-and-automation/a-practical-guide-to-getting-started-with-policy-as-code/>
- Golightly L, Modesti P, Garcia R, Chang V. Securing distributed systems: a survey on access control techniques for cloud, blockchain, IoT and SDN. *Cyber Secur Appl.* 2023;1:100015. doi:10.1016/j.csa.2023.100015
- Agarwal V, Butler C, Degenaro L, Kumar A, Sailer A, Steinder G. Compliance-as-code for cybersecurity automation in hybrid cloud. In: *Proc IEEE Int Conf Cloud Comput (CLOUD)*; 2022. p. 427–437. doi:10.1109/CLOUD55607.2022.00066
- Luo T, Pan M, Tholoniati P, Cidon A, Geambasu R, Lécuyer M. Privacy budget scheduling. In: *Proc 15th USENIX Symp Operating Systems Design and Implementation (OSDI 21)*; 2021. Available from: <https://www.usenix.org/system/files/osdi21-luo.pdf>
- National Institute of Standards and Technology (NIST). Adversarial machine learning: a taxonomy and terminology of attacks and mitigations. NIST AI 600-1; 2024. Available from: <https://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.600-1.pdf>
- Chen Y, Zhao Y, Li X, Zhang J, Long J, Zhou F. An open dataset of data lineage graphs for data governance research. *Vis Inform.* 2024;8(1):1–5. doi:10.1016/j.visinf.2024.01.001
- Pittala SK. Cybersecurity and online safety: a critical asset in the information era. *J Front Multidiscip Res.* 2023;4(1):576. doi:10.54660/jfmr.2023.4.1.576-579
- Confidential Computing Consortium. Confidential computing: hardware-based trusted execution for applications and data. Whitepaper; 2022 Nov (updated). Available from: https://confidentialcomputing.io/wp-content/uploads/sites/10/2023/03/CCC_outreach_whitepaper_updated_November_2022.pdf
- Ray J. *Policy as code: improving cloud native security.* O'Reilly Media; 2024. ISBN: 978-1098139186.
- Ashok VKC. Cybersecurity for smart infrastructure and public utilities. *Int J Multidiscip Res Growth Eval.* 2023;4(2):947–949. doi:10.54660/IJMRGE.2023.4.2.947-949
- Santos Silva R, Brito A, de Lima Filho JP. Automation of security controls for continuous compliance in vulnerability management. In: *Proc LADC*; 2024. doi:10.1145/3697090.3697107
- Pittala SK, Ashok VKC. A new era in security: bridging information security and cybersecurity. *Int J Multidiscip Futur Dev.* 2023;4(1):69–72. doi:10.54660/IJMF.2023.4.1.69-72
- Gunda SK. Fault prediction unveiled: analyzing the effectiveness of random forest, logistic regression, and K-neighbors. In: *Proc 2nd Int Conf Self Sustainable Artificial Intelligence Systems (ICSSAS)*; 2024; Erode, India. p. 107–113. doi:10.1109/ICSSAS64001.2024.10760620
- Sivva SD, Thalakanti RR, Bandari SSG, Yettapu SDR. AI-driven decision intelligence for agile software

- lifecycle governance: an architecture-centered framework integrating machine learning defect prediction and automated testing. *Int J Emerg Technol Comput Sci Inf Technol*. 2023;4(4):167–172. Available from: <https://www.ijetcsit.org/index.php/ijetcsit/article/view/554>
27. Shimamura S. A story of introducing data lineage into LINE's large-scale data platform. LINE Engineering Blog. 2022 Nov 24. Available from: <https://engineering.linecorp.com/en/blog/data-lineage-on-line-big-data-platform/>
28. Cutler JW, Disselkoe C, Eline A, He S, Headley K, Hicks M, *et al*. Cedar: a new language for expressive, fast, safe, and analyzable authorization (extended version). *arXiv*. 2024. doi:10.48550/arXiv.2403.04651. Available from: <https://arxiv.org/abs/2403.04651>