

Original Article

AI Governance by Design: Embedding Trust, Compliance, and Auditability into Enterprise AI Platforms

Hemant Soni

Principal Telecom Solution Architect, Capgemini America Inc., Atlanta, Georgia, USA

Received Date: 14 June 2026

Revised Date: 19 June 2026

Accepted Date: 29 June 2026

Abstract: Enterprises across regulated industries now build and operate artificial intelligence on shared platforms that host dozens or hundreds of models. Governance, the work of keeping those systems trustworthy, compliant, and answerable to scrutiny, is too often added after the platform is built, when it is expensive and sometimes impossible to retrofit. This paper proposes a governance-by-design framework for enterprise AI platforms that treats trust, compliance, and auditability as platform properties rather than late-stage controls. We synthesize guidance from the NIST AI Risk Management Framework, the European Union AI Act, ISO/IEC 42001, the OECD AI Principles, and the data protection rules of the GDPR, then map that guidance onto a four-layer reference architecture and a lifecycle loop that produces audit evidence as the system runs. The framework specifies where policy enforcement, risk tiering, model documentation, and evidence capture attach, and it draws on internal algorithmic auditing to make verification routine rather than exceptional. A qualitative comparison with retrofitted governance indicates that designing for governance early gives broader coverage across six evaluation dimensions and, in particular, makes auditability feasible at lower cost. We discuss limitations, including the gap between modeled and measured assurance and the overhead of evidence capture. The contribution is a standards-aligned blueprint that enterprise architects can adapt for regulated, high-volume AI.

Keywords: AI Governance, Governance by Design, Auditability, Regulatory Compliance, Enterprise AI Platforms, NIST AI RMF.

I. INTRODUCTION

Enterprises have industrialized AI. A modern organization no longer runs a single model; it runs a portfolio of them on shared platforms built from model registries, feature stores, and automated pipelines. In sectors such as telecom and financial services, hundreds of models may decide credit, route customers, flag fraud, and tune networks at the same time. The platform, rather than any one model, has become the unit that must be governed. Governance is the work of keeping that portfolio trustworthy, compliant, and answerable. In practice it is often the last thing added. Teams build a model, prove it performs, ship it, and only then ask how to document it, how to show a regulator it is lawful, and how to reconstruct what it did when something goes wrong. By then the cheapest moment had passed. Auditability is the clearest example, since an audit trail that was never captured cannot be recreated after this fact.

Regulation has raised the cost of getting this wrong. The European Union AI Act sets tiered obligations, with heavy documentation and oversight duties for high-risk uses [5]. The GDPR already requires data protection by design [8]. Sector regulators add rules of their own. An enterprise platform that serves many business lines can fall under several regimes at once, so governance bolted on model by model multiplies effort instead of sharing it.

The literature on what good AI should look like is rich. Principle sets from the European expert group [1], from Floridi and colleagues [2], and from the OECD [3] converge on a familiar list: human oversight, robustness, privacy, transparency, fairness, and accountability. Gasser and Almeida proposed a layered model that separates the technical, ethical, and social levels at which AI must be governed [10]. On the practical side, documentation tools such as model cards [12] and datasheets for datasets [13] make systems legible, internal algorithmic auditing defines how to verify them before and after release [11], and established methods address fairness [14], explanation [15], adversarial robustness [16], and privacy, including differential privacy [17] and federated learning [18]. Standards have begun to consolidate the field. The NIST AI Risk Management Framework organizes the work into govern, map, measure, and manage [4]; ISO/IEC 42001 defines an auditable AI management system [6] that sits alongside the established security standard ISO/IEC 27001 [7]; and privacy by design is now both a principle and a legal duty [8], [9]. Telecom bodies add domain guidance through the ETSI work on securing AI [19] and the GSMA responsible AI roadmap [20].



What is missing is a treatment that binds this guidance to the platform layer where enterprises actually build, and that treats auditability as a first-class property rather than a reporting chore. Most frameworks are written for AI in the abstract or for a single model. The aim of this paper is to create a governance-by-design framework for enterprise AI platforms, organized around three pillars: trust, compliance, and auditability. Our hypothesis is that embedding these pillars in the platform architecture and lifecycle improves coverage across recognized evaluation dimensions and, above all, makes credible auditability achievable at a cost that retrofitting cannot match. Telecom, the author's own domain, serves as a running example.

II. METHODOLOGY

This is a conceptual contribution, so the method is a design-science-informed synthesis rather than an experiment. It proceeds in five steps: select authoritative sources, favoring standards and peer-reviewed work; extract the governance pillars and their controls; derive a platform reference architecture and a lifecycle loop; map the pillars to design mechanisms and to the evidence they produce; and define evaluation criteria to compare governance by design against retrofitted governance. We organize governance around three pillars. Trust is the property that a system behaves as intended and that stakeholders can rely on it, with robustness, fairness, transparency, security, and privacy rolling up into it. Compliance is the property that the system demonstrably satisfies the laws, regulations, and standards that apply to it. Auditability is the property that the system continuously produces the evidence needed to verify trust and compliance, independently and after the fact. The order is deliberate. Auditability is what turns the other two claims into something a regulator, a customer, or an internal reviewer can actually check.

Figure 1 organizes an enterprise AI platform into four layers. A data foundation governs how information enters the system and is tracked. A platform layer hosts model development and deployment through MLOps. A governance and control plane holds the policy, risk, and evidence machinery. An application layer delivers the AI services that the business ships. Trust, compliance, and auditability run as cross-cutting concerns through all four layers. The control plane is the new emphasis: it is where policy as code, risk tiering, approval gates, and audit logging live, so governance is enforced by the platform rather than left to each team.



Figure 1: Four-Layer Reference Architecture for an Enterprise AI Platform, With Trust, Compliance, and Auditability as Cross-Cutting Design Principles

Figure 2 recasts the same controls as a lifecycle loop that generates evidence at every turn. Framing sets the purpose and the risk tier and writes the control plan. The data stage records consent, lineage, and quality. The build stage applies controls, testing, and registry entries. Deployment passes through approval gates and publishes model documentation. Operation monitors behavior and emits audit evidence that flows back into the next cycle. The central idea is that evidence is a by-product of running the system, captured as it happens, not assembled under deadline once an audit is announced.



Figure 2: Governance Lifecycle Loop in Which Each Stage Applies Controls and Emits Audit Evidence That Feeds the Next Cycle

Table 1 ties the three pillars to the design mechanism that delivers each and, importantly, to the audit evidence that mechanism produces. The evidence column is what makes the framework auditable, because each control leaves a trace that an independent reviewer can inspect.

Table 1: Governance Pillars Mapped To By-Design Mechanisms and the Audit Evidence They Produce

Pillar	By-design mechanism	Audit evidence produced
Trust	Robustness and fairness testing, human oversight	Test reports, oversight and override logs
Compliance	Risk tiering, policy as code, control mapping	Conformity records, policy decision logs
Auditability	Lineage capture, model registry, immutable logs	Versioned data and model lineage, audit trail

III. RESULTS

The synthesis yields three results: the layered architecture and lifecycle, a comparison of governance by design against retrofitted governance, and a mapping to standards and regulation. These results are analytical and should be read as structured reasoning from the source guidance rather than field measurements.

Architectural result Placing governance in the control plane gives every pillar a home in the platform itself. Policy as code turns written rules into checks that run on each pipeline. Risk tiering decides up front how much scrutiny a use case need. Lineage and immutable logging mean the evidence trail already exists by the time anyone asks for it. Because the control plane sits beside the model registry, one governed pipeline can serve many models, so the marginal cost of governing the next model falls instead of repeating per project.

Comparative result Figure 3 contrasts governance by design with a retrofitted approach across six dimensions on an illustrative five-point scale. The by-design profile is fuller and more even. The retrofitted profile sags most on auditability, the dimension that depends most on capturing evidence while the system runs. Compliance and trust also lag, since both rest on decisions such as risk tiering and data lineage that are hard to reconstruct later. Robustness and efficiency narrow the gap a little, because some testing and tuning can be added after the fact.

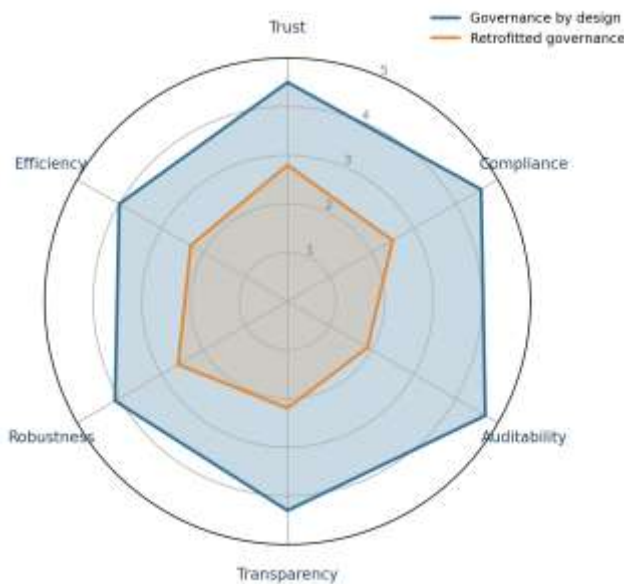


Figure 3: Illustrative Comparison of Coverage across Six Dimensions for Governance by Design versus Retrofitted Governance

Table 2 states the same comparison in words, naming the mechanism a by-design platform uses for each dimension and the typical weakness of retrofitting. The pattern is not that retrofitting fails outright, but that it delivers partial coverage at higher cost, with auditability the hardest to recover.

Table 2: Evaluation Dimensions: Governance-By-Design Mechanism and the Common Limitation of a Retrofit

Dimension	Governance-by-design mechanism	Retrofit limitation
Trust	Testing and oversight built into pipelines	Behavior accepted without a baseline
Compliance	Policy as code and risk tiering up front	Manual mapping after the fact
Auditability	Evidence captured as the system runs	Trail cannot be reconstructed
Transparency	Documentation generated at each stage	Rationale lost over time
Robustness	Adversarial and drift testing in build	Reactive patching only
Efficiency	Shared control plane across models	Repeated per-project effort

The third result anchors the framework to external requirements. Table 3 maps the four platform layers to the standards and regulations that govern them. For an enterprise platform the mapping is practical, because a single well-designed control plane can satisfy several audits at once instead of meeting each regime separately.

Table 3: Mapping Platform Layers to Standards and Regulatory Anchors

Platform layer	Primary anchors	Governance focus
Data foundation	GDPR, ISO/IEC 27001	Lawful, traceable, secure data
AI and ML platform	NIST AI RMF (Measure), ETSI SAI	Tested, secure model building
Governance and control of plane	ISO/IEC 42001, EU AI Act	Risk tiering, oversight, audit evidence
Applications and services	EU AI Act, OECD AI Principles	Acceptable use and disclosure

IV. DISCUSSION

The analysis supports the hypothesis in a qualitative sense. Governance by design gives broader and more even coverage, and the advantage is largest exactly where it matters most for regulated enterprises, namely auditability. The reason is structural. Audit evidence is a record of events, and if the platform did not capture it as those events happened, no later effort recovers it faithfully. The same logic applies to data lineage and to the risk decisions on which compliance rests.

For enterprise architects the practical implication is that the governance machinery belongs in the platform, not in a separate compliance function that reviews finished systems. A control plane with policy as code, a model registry, risk tiering, and immutable logging changes what teams can and cannot do by default. The standards mapping in Table 3 suggests this is also efficient: build the control plane once and it answers to several regimes. For a telecom operator that already lives under sector regulation and now faces AI rules, that shared layer is the difference between governance that scales and governance repeated, expensively, for every model.

Against prior research, the contribution is integrative and platform-specific rather than a new technique. The principles [1], [2], [3], the layered governance model [10], the documentation and auditing practices [11], [12], [13], and the privacy and robustness methods [16], [17], [18] are all established. What the framework adds is their placement inside an enterprise platform, with auditability treated as a first-class pillar. General frameworks such as the NIST AI RMF [4] are deliberately domain-neutral; this paper is one attempt to operationalize that guidance at the platform layer.

The limitations are real. The comparison in Figure 3 is illustrative, not measured, and a convincing test would instrument a live platform and report evidence coverage and audit outcomes. Capturing evidence carries a cost in storage, performance, and engineering effort, and over-instrumentation can create risk of its own by collecting more than is needed. The standards landscape is still moving, so any fixed mapping will age. Agreed metrics are scarce, especially for auditability and fairness. And the framework assumes a level of platform maturity, including an MLOps foundation and a willingness to enforce policy centrally, that not every enterprise has.

V. CONCLUSIONS

Governance is decided when a platform is designed, not when an audit arrives. This paper has argued that enterprises should treat trust, compliance, and auditability as platform properties, embedded in a four-layer architecture and a lifecycle loop that produces evidence as a matter of course. The framework consolidates established guidance and places it in the control plane where enterprise teams build, with a mapping to the standards and regulations that increasingly govern these systems. The

qualitative analysis indicates that designing for governance early gives broader coverage and, above all, makes credible auditability achievable at a cost that retrofitting cannot match.

Two lines of future work follow. The first is empirical validation through a pilot that instruments a live platform and measures evidence coverage, audit effort, and outcomes rather than modeled scores. The second is a reference implementation of the control plane, built on policy as code, together with metrics for auditability that turn the pillar from a principle into a measurable practice. With those in place, the blueprint offered here can help regulated, high-volume enterprises adopt AI they can stand behind.

VI. ACKNOWLEDGEMENTS

This work received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors. The author thanks colleagues in the enterprise architecture and telecom community for discussions that shaped the treatment of governance as a design concern, and the reviewers whose comments improved the clarity of the argument. The views expressed are those of the author and do not necessarily represent those of any affiliated organization.

VII. REFERENCES

- [1] European Commission High-Level Expert Group on Artificial Intelligence, "Ethics Guidelines for Trustworthy AI," Brussels, 2019.
- [2] L. Floridi, J. Cows, M. Beltrametti, R. Chatila, et al., "AI4People: An Ethical Framework for a Good AI Society," *Minds and Machines*, vol. 28, no. 4, pp. 689-707, 2018.
- [3] Organization for Economic Co-operation and Development, "Recommendation of the Council on Artificial Intelligence," OECD/LEGAL/0449, Paris, 2019.
- [4] National Institute of Standards and Technology, "Artificial Intelligence Risk Management Framework (AI RMF 1.0)," NIST AI 100-1, Gaithersburg, MD, 2023.
- [5] European Parliament and Council, "Regulation (EU) 2024/1689 Laying Down Harmonised Rules on Artificial Intelligence (AI Act)," *Official Journal of the European Union*, 2024.
- [6] International Organization for Standardization, "ISO/IEC 42001:2023, Information Technology, Artificial Intelligence, Management System," Geneva, 2023.
- [7] International Organization for Standardization, "ISO/IEC 27001:2022, Information Security, Cybersecurity and Privacy Protection, Information Security Management Systems," Geneva, 2022.
- [8] European Parliament and Council, "Regulation (EU) 2016/679, General Data Protection Regulation," *Official Journal of the European Union*, 2016.
- [9] A. Cavoukian, "Privacy by Design: The 7 Foundational Principles," Information and Privacy Commissioner of Ontario, Toronto, 2011.
- [10] U. Gasser and V. A. F. Almeida, "A Layered Model for AI Governance," *IEEE Internet Computing*, vol. 21, no. 6, pp. 58-62, 2017.
- [11] I. D. Raji, A. Smart, R. N. White, M. Mitchell, et al., "Closing the AI Accountability Gap: Defining an End-to-End Framework for Internal Algorithmic Auditing," in *Proc. ACM Conf. on Fairness, Accountability, and Transparency*, 2020, pp. 33-44.
- [12] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, et al., "Model Cards for Model Reporting," in *Proc. ACM Conf. on Fairness, Accountability, and Transparency*, 2019, pp. 220-229.
- [13] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, et al., "Datasheets for Datasets," *Communications of the ACM*, vol. 64, no. 12, pp. 86-92, 2021.
- [14] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, "A Survey on Bias and Fairness in Machine Learning," *ACM Computing Surveys*, vol. 54, no. 6, pp. 1-35, 2021.
- [15] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why Should I Trust You? Explaining the Predictions of Any Classifier," in *Proc. 22nd ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, 2016, pp. 1135-1144.
- [16] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and Harnessing Adversarial Examples," in *Proc. Int. Conf. on Learning Representations (ICLR)*, 2015.
- [17] C. Dwork and A. Roth, "The Algorithmic Foundations of Differential Privacy," *Foundations and Trends in Theoretical Computer Science*, vol. 9, no. 3-4, pp. 211-407, 2014.
- [18] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. Aguera y Arcas, "Communication-Efficient Learning of Deep Networks from Decentralized Data," in *Proc. AISTATS*, 2017.
- [19] ETSI Industry Specification Group on Securing Artificial Intelligence, "Securing Artificial Intelligence (SAI); Problem Statement," ETSI GR SAI 004 V1.1.1, Sophia Antipolis, 2020.
- [20] GSMA, "The GSMA Responsible AI Maturity Roadmap," London, 2024.