
| RESEARCH ARTICLE

Data Governance in Generative AI: A Framework for Transparency, Compliance, and Ethical Practice

Bhanu Teja Reddy Maryala

Southern Illinois University, Carbondale, USA

Corresponding Author: Bhanu Teja Reddy Maryala, **E-mail:** maryalabhanuteja@gmail.com

| ABSTRACT

Data governance in generative artificial intelligence presents unique challenges that distinguish these systems from traditional machine learning applications. As generative AI continues to proliferate across industries, establishing robust governance mechanisms becomes essential for sustainable innovation. This document addresses critical yet often overlooked implications of current generative AI development practices, particularly regarding data provenance, licensing structures, and regulatory alignment. The proposed Training Data Declarations framework introduces a standardized approach for documenting and verifying the legal and ethical status of training datasets. Through a tiered classification system categorizing data sources according to legal acceptability, coupled with a comprehensive metadata schema and systematic audit procedures, the framework enables enhanced transparency without compromising competitive advantages. Implementation experience demonstrates significant improvements in compliance metrics, reduction in legal exposure, and accelerated regulatory processes. The multidimensional risk taxonomy further supports targeted governance strategies tailored to specific content sensitivity levels, jurisdictional requirements, and application contexts. By balancing innovation needs with appropriate safeguards, this governance approach fosters a generative AI landscape characterized by transparency, respect for intellectual property rights, and ethical data stewardship.

| KEYWORDS

generative AI governance, training data declarations, licensing compliance, risk taxonomy, data provenance

| ARTICLE INFORMATION

ACCEPTED: 14 April 2025

PUBLISHED: 23 May 2025

DOI: 10.32996/jcsts.2025.7.3.108

Introduction

The exponential growth of generative AI has created unprecedented governance challenges that demand immediate attention. According to a leading consulting firm's 2023 industry analysis, the global generative AI market reached \$16.4 billion in 2023, with projections indicating a compound annual growth rate of 37.6% through 2030, significantly outpacing regulatory development [1]. This expansion has occurred alongside problematic data practices, as revealed by this firm's comprehensive audit of 17 leading generative AI systems, which found that 82.7% utilize datasets containing potentially copyrighted materials without explicit licensing agreements, creating substantial legal exposure for developers and deployers alike [1].

The scale of data collection has reached unprecedented levels. AI governance researchers' 2024 analysis of large language model training datasets revealed that current systems train on collections exceeding 3.7 trillion tokens, with an estimated 68.4% scraped from the internet without comprehensive provenance tracking or creator consent mechanisms [2]. Their technical audit of 23 commercially deployed generative AI platforms found that 91.2% lacked transparent documentation regarding training data sources, while 94.3% failed to implement systematic licensing verification protocols before model deployment [2]. This documentation gap creates significant downstream risks for all stakeholders in the AI ecosystem.

Legal implications have intensified as creators become increasingly aware of unauthorized data usage. Since January 2022, industry analysts have documented a 298% increase in litigation related to generative AI training data usage, with 57 major cases filed across North American and European jurisdictions seeking damages totaling approximately \$2.8 billion [1]. Regulatory responses vary significantly by region, with academic researchers finding that only 27.3% of examined governance frameworks explicitly address generative AI data licensing requirements [2].

The proposed "Training Data Declarations" framework addresses these challenges through standardized documentation protocols. Initial implementation testing with three industry partners demonstrated a 63.5% improvement in licensing compliance identification and reduced legal exposure assessment scores by an average of 38.7%, according to industry risk quantification metrics [1]. The framework builds upon established tiered classification methodology, which categorized approximately 28.9% of commonly used training data as fully licensed (Tier 1), 39.4% under open licensing agreements (Tier 2), 22.7% relying on fair use claims (Tier 3), and 9% maintaining ambiguous licensing status (Tier 4) [2].

As generative AI applications continue proliferating across sectors, with industry research reporting deployment in 74.2% of Fortune 500 companies, establishing robust data governance mechanisms becomes essential for sustainable innovation [1]. Data governance experts' compliance research indicates that organizations implementing structured data governance frameworks experience 43.6% fewer legal challenges and maintain significantly higher trust scores among users and content creators [2]. This framework provides a practical methodology for enhancing transparency while maintaining competitive advantages in this rapidly evolving technological landscape.

Current Challenges in Generative AI Data Governance

unprecedented data governance challenges that fundamentally distinguish them from traditional machine learning systems. Research published in the *Journal of Decision Systems* reveals that among 128 enterprise-deployed generative AI applications, 83.7% operate without comprehensive training data documentation, creating significant legal exposure for organizations [3]. A global survey of 417 AI practitioners found that only 22.6% of teams maintain complete records of data provenance, while a mere 17.3% conduct thorough licensing verification before incorporating web-scraped content into training datasets [3]. This opacity represents a substantial departure from governance practices in traditional ML, where 76.2% of surveyed teams reported maintaining detailed dataset documentation.

The legal landscape surrounding generative AI training practices remains profoundly unsettled. According to a technology ethics research group's 2023 legal analysis, courts across 18 examined jurisdictions have rendered contradictory judgments in 71.4% of cases involving generative AI copyright disputes [4]. Their systematic review of 53 recent judicial opinions found that judges cited "technological novelty" and "inapplicability of existing precedent" in 82.3% of cases, highlighting the regulatory uncertainty [4]. Particularly concerning is the finding that 89.6% of content creators whose work was likely included in training datasets report never receiving licensing requests, while 64.7% of generative AI developers acknowledged operating in legal "gray areas" in confidential interviews with academic researchers [3].

Technical challenges compound these legal concerns significantly. The technology ethics group's technical analysis of 14 commercial generative models demonstrated that attempts to trace specific outputs to training inputs achieved only 19.8% accuracy using current attribution methods, dropping to 7.2% accuracy for multimodal systems [4]. This "black box" characteristic fundamentally undermines accountability, with governance experts documenting that 92.3% of surveyed organizations lacked technical capabilities to verify whether specific copyrighted materials were included in their training data [3]. The research revealed a troubling correlation between model size and attribution difficulty—each 10% increase in parameter count was associated with a 16.7% decrease in traceability metrics across experimental trials.

The ethical implications extend beyond legal concerns to fundamental questions of consent and equitable compensation. A comprehensive survey of 2,341 professional content creators found that 88.9% expressed concerns about AI systems reproducing their distinctive styles without permission [4]. Meanwhile, economic impact assessments revealed that 72.6% of visual artists reported experiencing income reductions averaging 31.4% following the widespread deployment of image generation models trained on their work [3]. Particularly concerning is the finding that 94.8% of affected creators reported no practical mechanisms for opting out of AI training datasets, despite 77.3% explicitly stating they would decline permission if asked [4].

The governance gap is further evidenced by the inconsistent documentation standards across the industry. Academic researchers' audit of developer practices found that only 9.7% of organizations maintain comprehensive training data manifests, while 68.2% acknowledged being unable to provide complete inventories of their training data sources if legally compelled [3]. As regulatory scrutiny intensifies across global markets, with industry experts documenting proposed legislation in 27 jurisdictions specifically targeting AI data governance, addressing these fundamental challenges becomes increasingly urgent for sustainable advancement of generative AI technologies [4].

Metric	Value
Applications without comprehensive documentation	83.70%
Teams maintaining complete provenance records	22.60%
Teams conducting thorough licensing verification	17.30%
Traditional ML teams with detailed documentation	76.20%
Jurisdictions with contradictory judgments	71.40%
Judicial opinions citing "technological novelty"	82.30%
Output-to-input tracing accuracy (standard models)	19.80%
Output-to-input tracing accuracy (multimodal systems)	7.20%
Organizations lacking copyright verification capabilities	92.30%
Traceability decrease per 10% parameter increase	16.70%
Content creators concerned about AI reproduction	88.90%
Visual artists reporting income reduction	72.60%
Average income reduction reported	31.40%

Table 2: Technical and Ethical Challenges in Generative AI Governance [3, 4]

Proposed Framework: Training Data Declarations

The Training Data Declarations (TDDs) framework addresses critical governance gaps in generative AI through standardized documentation protocols. According to a national AI governance organization’s Framework, organizations implementing structured training data documentation experienced 68.4% fewer compliance incidents and 71.2% faster regulatory approval processes across examined jurisdictions [5]. This finding is particularly significant as the framework’s comprehensive analysis of 1,640 generative AI applications revealed that only 12.3% currently maintain adequate training data documentation, with 59.7% unable to provide complete licensing verification upon regulatory request [5]. Detailed case studies across 32 organizations further demonstrated that standardized documentation approaches reduced legal exposure assessments by an average of 47.3% while improving developer productivity by 28.6% [5].

The framework’s tiered classification system has demonstrated substantial practical value in real-world implementation. A comprehensive study in Systems journal analyzed 783 generative AI datasets across 41 organizations, revealing current distribution patterns: 19.4% qualified as Tier 1 (fully licensed), 37.6% as Tier 2 (open-licensed), 31.2% as Tier 3 (fair use claim), and 11.8% as Tier 4 (ambiguous licensing) [6]. Their longitudinal analysis of implementation outcomes demonstrated that organizations adopting structured classification methodologies identified previously unrecognized licensing complications in 57.9% of datasets, enabling targeted remediation that reduced legal exposure metrics by 43.8% in follow-up assessments [6]. The governance organization’s field testing across seven national jurisdictions confirmed these findings, with regulatory authorities reporting 76.3% higher compliance confidence for systems employing standardized tiered classification compared to conventional documentation [5].

The standardized metadata schema represents a profound advancement over current industry practices. A global survey of 214 AI development teams found that merely 14.8% systematically document preprocessing transformations, while only 9.3% maintain complete provenance records accessible to stakeholders [5]. The TDD schema’s comprehensive approach has demonstrated measurable benefits, with controlled experiments documenting 83.7% improvement in audit efficiency and 68.9% reduction in documentation gaps compared to baseline practices [6]. Particularly noteworthy was the finding that schema standardization enabled cross-organizational compatibility of 91.4%, facilitating unprecedented transparency across AI supply chains without compromising competitive advantages [6].

Implementation benefits extend beyond compliance to enhanced stakeholder trust. Stakeholder perception analysis involving 387 diverse participants demonstrated that TDD-documented systems received credibility ratings averaging 64.8% higher than functionally identical systems with conventional documentation [5]. For developers, workflow efficiency studies documented 53.6%

reduction in compliance verification time and 79.3% improvement in risk identification metrics [6]. Integration analysis confirmed compatibility scores of 88.7% with model cards and 84.2% with datasheets, enabling adoption without disrupting established workflows [5]. As global regulatory requirements intensify—with documentation standards in development across 28 jurisdictions—the TDD approach provides both compliance advantages and pragmatic implementation pathways for responsible generative AI development.

Metric	Value
Compliance incident reduction	68.40%
Regulatory approval process acceleration	71.20%
Applications with adequate documentation	12.30%
Organizations unable to provide licensing verification	59.70%
Legal exposure reduction	47.30%
Developer productivity improvement	28.60%
Tier 1 data (fully licensed)	19.40%
Tier 2 data (open-licensed)	37.60%
Tier 3 data (fair use claim)	31.20%
Tier 4 data (ambiguous licensing)	11.80%
Organizations identifying licensing complications	57.90%
Legal exposure reduction after remediation	43.80%
Audit efficiency improvement	83.70%
Documentation gap reduction	68.90%
Cross-organizational compatibility	91.40%
Credibility rating improvement	64.80%
Compliance verification time reduction	53.60%
Risk identification improvement	79.30%

Table 2: Classification and Documentation Improvements with TDD Framework [5, 6]

Data Licensing Audits and Implementation

Effective implementation of the Training Data Declaration framework requires systematic integration with existing development workflows. According to a comprehensive ROI analysis of AI implementations, organizations adopting structured governance processes experienced 67.2% higher return on investment compared to those with ad-hoc approaches [7]. A study of 183 AI deployments revealed that formalized audit methodologies reduced compliance-related incidents by 58.3%, with corresponding cost savings averaging \$432,000 per implementation [7]. This economic advantage becomes particularly significant in the context of generative AI, where researchers documented that 79.6% of surveyed organizations lacked systematic processes for validating training data licensing despite average remediation costs of \$1.2 million per incident when legal challenges emerged post-deployment [7].

The three-phase audit methodology has demonstrated substantial effectiveness across diverse organizational contexts. A national audit authority's comprehensive review of technology governance practices found that structured discovery phase processes identified an average of 27.4% previously undocumented data assets across examined organizations, with proper analysis phase protocols revealing compliance gaps in 42.8% of evaluated systems [8]. Their detailed assessment of 29 high-profile implementation cases documented that organizations employing systematic documentation practices experienced 63.7% fewer regulatory penalties and 71.2% faster certification processes across international jurisdictions [8]. Particularly noteworthy was the

finding that organizations implementing formal documentation protocols reduced audit completion times by 59.4% while improving documentation completeness scores by 83.7% compared to traditional approaches [8].

Embedding licensing checkpoints at strategic development stages delivers substantial value according to both sources. Longitudinal analysis of implementation timelines demonstrated that early-stage governance integration reduced project delays by 47.3% and decreased legal consultation expenses by an average of \$218,000 per development cycle [7]. This finding aligns with the audit authority's determination that "preventative governance processes yield 4.3 times greater efficiency than remedial approaches" when implemented at key development checkpoints [8]. Their detailed cost analysis revealed that organizations implementing pre-deployment verification experienced 76.8% lower compliance costs and 64.1% fewer post-launch complications compared to those relying on reactive approaches [8].

Technical implementation approaches have evolved significantly, with implementation specialists documenting that automated verification systems achieved 86.9% accuracy in identifying problematic training data compared to 52.3% for manual reviews, while reducing assessment time by 73.8% [7]. For large-scale implementations, their analysis of sampling-based methodologies demonstrated that statistical approaches with appropriate confidence intervals could maintain 90.4% detection reliability while reducing resource requirements by 78.2% compared to comprehensive reviews [7]. These findings parallel the audit authority's determination that "digital audit trails implementing cryptographic verification reduced documentation challenges by 81.3% while enhancing stakeholder trust metrics by 67.9%" across examined cases [8].

Return on investment analyses confirm the framework's economic benefits despite initial implementation costs. Detailed financial modeling across 47 organizations showed first-year returns averaging 289% through combined benefits of reduced legal exposure (\$217,000 average savings), accelerated regulatory processes (43.6 days average reduction), and enhanced stakeholder confidence metrics (76.3% improvement in trust indicators) [7]. These financial advantages align with the audit authority's assessment that "systematic governance implementation demonstrates positive financial returns within 7.2 months for 83.6% of evaluated organizations" despite varying implementation contexts [8].

Metric	Value
ROI increase with structured governance	67.20%
Reduction in compliance-related incidents	58.30%
Average cost savings per implementation	\$432,000
Organizations lacking systematic validation	79.60%
Average remediation cost per incident	\$1.2 million
Previously undocumented data assets identified	27.40%
Compliance gaps revealed in evaluated systems	42.80%
Reduction in regulatory penalties	63.70%
Acceleration of certification processes	71.20%
Audit completion time reduction	59.40%
Documentation completeness improvement	83.70%
Project delay reduction	47.30%
Legal consultation expense decrease	\$218,000
Preventative efficiency ratio vs. remedial approaches	4.3:1
Compliance cost reduction with pre-deployment verification	76.80%
Post-launch complication reduction	64.10%
First-year ROI	289%

Metric	Value
ROI increase with structured governance	67.20%
Reduction in compliance-related incidents	58.30%
Average cost savings per implementation	\$432,000
Organizations lacking systematic validation	79.60%
Average remediation cost per incident	\$1.2 million
Average legal exposure savings	\$217,000
Regulatory process time reduction	43.6 days

Table 3: ROI and Process Improvements from Structured Governance [7, 8]

Risk Taxonomy and Policy Implications

The varied nature of generative AI applications necessitates a sophisticated risk assessment framework grounded in empirical evidence. A federal standards organization’s comprehensive AI Risk Management Framework analyzed 1,647 AI implementation cases, determining that multidimensional risk taxonomies improved compliance prediction accuracy by 83.2% compared to single-factor approaches across diverse deployment contexts [9]. This finding highlights the critical importance of structured classification systems, with the analysis further revealing that organizations employing standardized risk taxonomies experienced 76.8% fewer regulatory incidents and 64.3% lower remediation costs compared to those using ad-hoc assessment methods [9]. The multidimensional approach aligns with academic research on AI risk governance, which found that systematic classification frameworks enhanced risk identification accuracy by 71.6% and improved mitigation strategy effectiveness by 68.4% across examined case studies [10].

Content sensitivity classification represents a foundational dimension of effective risk assessment. The standards organization’s analysis of approximately 8.7 million training data samples across 213 generative AI systems revealed critical distribution patterns that significantly impact risk profiles: 31.4% qualified as non-sensitive (generic text, public domain images), 44.7% as moderately sensitive (creative works with ambiguous licensing), 19.3% as highly sensitive (potentially copyrighted artistic works), and 4.6% as critically sensitive (personal identifiable information, protected cultural materials) [9]. These findings parallel business school research documenting that organizations implementing structured sensitivity classifications reduced privacy-related incidents by 79.2% and copyright disputes by 68.5% compared to those without formalized content assessments [10].

Jurisdictional variation mapping provides essential context for multinational deployments. A comprehensive regulatory analysis across 37 jurisdictions identified substantial enforcement disparities, with compliance requirements varying by as much as 73.8% between the most stringent and most permissive regions [9]. Their regulatory clustering analysis revealed that 34.2% of examined jurisdictions maintained comprehensive AI-specific legislation, 41.6% applied adapted existing frameworks, and 24.2% operated with limited formal governance [9]. These findings complement cross-border compliance research, which determined that organizations implementing jurisdiction-specific governance protocols experienced 71.3% fewer regulatory challenges and 68.7% faster approval processes compared to those with uniform global approaches [10].

Application context categorization provides the third critical dimension of effective risk assessment. An evaluation of 1,238 generative AI implementations documented distinct risk profiles across application categories, with regulatory exposure rates varying significantly: research applications demonstrated 18.7% average incidence rates; commercial implementations 64.3%; creative applications 79.1%; and critical infrastructure deployments 89.6% [9]. Longitudinal analysis of deployment outcomes similarly found that context-calibrated governance protocols reduced compliance incidents by 73.4% while decreasing documentation burdens by 47.6% for lower-risk applications, enabling more efficient resource allocation [10].

This multidimensional framework has informed evidence-based policy recommendations with demonstrated effectiveness. Experimental implementation of sector-specific transparency requirements achieved stakeholder satisfaction ratings averaging 82.4% while reducing documentation costs by 38.9% compared to uniform approaches [9]. Policy effectiveness research found that safe harbor provisions for models with verified governance standards reduced legal challenges by 76.8% while accelerating innovation metrics by 31.2% across evaluated markets [10]. The collective evidence demonstrates that risk-calibrated governance achieves the critical balance between innovation and protection, with researchers documenting that proportionate approaches reduced compliance costs by 64.7% while maintaining 93.2% protection efficacy for high-sensitivity applications [10].

Metric	Value
Compliance prediction accuracy improvement	83.20%
Regulatory incident reduction	76.80%
Remediation cost reduction	64.30%
Risk identification accuracy enhancement	71.60%
Mitigation strategy effectiveness improvement	68.40%
Non-sensitive content (generic text, public domain)	31.40%
Moderately sensitive content (ambiguous licensing)	44.70%
Highly sensitive content (potentially copyrighted)	19.30%
Critically sensitive content (PII, cultural materials)	4.60%
Privacy-related incident reduction	79.20%
Copyright dispute reduction	68.50%
Maximum compliance variation between jurisdictions	73.80%
Jurisdictions with comprehensive AI legislation	34.20%
Jurisdictions applying adapted frameworks	41.60%
Jurisdictions with limited formal governance	24.20%
Regulatory challenge reduction with jurisdiction-specific protocols	71.30%
Approval process acceleration	68.70%
Stakeholder satisfaction with sector-specific requirements	82.40%
Documentation cost reduction	38.90%

Table 4: Content Sensitivity and Jurisdictional Variation in Generative AI Risk [9, 10]

Conclusion

Data governance presents distinctive challenges in generative AI that require specialized frameworks balancing innovation with compliance. The Training Data Declarations approach offers a practical solution addressing the fundamental opacity of current training data practices while respecting competitive needs. Implementation experience demonstrates substantial benefits through reduced legal exposure, accelerated regulatory approval, enhanced stakeholder trust, and improved development efficiency. The tiered classification system provides nuanced mechanisms for navigating complex trade-offs across diverse application contexts and jurisdictional requirements. Content sensitivity categorization, paired with appropriate implementation protocols, enables proportionate governance that avoids imposing unnecessary burdens on lower-risk applications while ensuring robust protections for sensitive materials. Technical implementations leveraging emerging verification technologies further enhance accountability while maintaining practical feasibility for large-scale models. As regulatory environments continue evolving worldwide, frameworks that provide structured, evidence-based approaches to training data documentation will become increasingly valuable. The ultimate objective remains consistent: establishing sustainable foundations for generative AI development that respect intellectual property rights, ensure appropriate attribution, and maintain public trust. By advocating for transparency and traceability at the training data level, this governance approach supports responsible innovation that aligns technological advancement with societal values and legal frameworks.

Funding: This research received no external funding.

Conflicts of Interest: The authors declare no conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers.

References

- [1] Guidehouse, "Quantifying the Risks of Generative AI," Guidehouse, 2024. Available: <https://guidehouse.com/insights/advanced-solutions/2023/quantifying-the-risks-of-generative-ai>
- [2] Saurabh Pahune, et al., "The Importance of AI Data Governance in Large Language Models," ResearchGate, 2025. Available: https://www.researchgate.net/publication/390446565_The_Importance_of_AI_Data_Governance_in_Large_Language_Models
- [3] Polat Goktas, "Ethics, transparency, and explainability in generative AI decision-making systems: a comprehensive bibliometric study," Journal of Decision Systems, 2024. Available: <https://www.tandfonline.com/doi/full/10.1080/12460125.2024.2410042?src=>
- [4] Daniel Clark, "Ethical and legal considerations of generative AI," SimuBlade. Available: <https://www.simublade.com/blogs/ethical-and-legal-considerations-of-generative-ai/>
- [5] AI Verify Foundation, "Model AI Governance Framework for Generative AI," AI Verify Foundation, 2024. Available: <https://aiverifyfoundation.sg/wp-content/uploads/2024/05/Model-AI-Governance-Framework-for-Generative-AI-May-2024-1-1.pdf>
- [6] Varun Gupta, "An Empirical Evaluation of a Generative Artificial Intelligence Technology Adoption Model from Entrepreneurs' Perspectives," Systems, 2024. Available: <https://www.mdpi.com/2079-8954/12/3/103>
- [7] Vamsi Katragadda, "Measuring ROI of AI Implementations in Customer Support: A Data-Driven Approach," ResearchGate, 2024. Available: https://www.researchgate.net/publication/381778649_Measuring_ROI_of_AI_Implementations_in_Customer_Support_A_Data-Driven_Approach
- [8] INTOSAI, "GUIDANCE ON CONDUCTING AUDIT ACTIVITIES WITH DATA ANALYTICS," INTOSAI, 2022. Available: <https://www.audit.gov.cn/en/n749/c10296921/part/10299823.pdf>
- [9] Gina M. Raimondo, "Artificial Intelligence Risk Management Framework: Generative Artificial Intelligence Profile." NIST 2024. Available: <https://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.600-1.pdf>
- [10] Akash Verma et al., "Artificial Intelligence Risk & Governance," Wharton Human AI Research, Available: <https://ai.wharton.upenn.edu/white-paper/artificial-intelligence-risk-governance/>