# Autonomous Identity Governance
## AI-Driven Access Decisions and Self-Healing Systems

*With Human vs. AI Benchmark Study (n=34 Matched Environments)*

Automation Impact from 95 Enterprise IGA Programmes

### Kieran Upadrasta
CISSP, CISM, CRISC, CCSP | MBA | BEng

27 Years' Cyber Security | Big 4 (Deloitte, PwC, EY, KPMG) | 21 Years Financial Services
Professor of Practice, Schiphol University | Honorary Senior Lecturer, Imperials | UCL Researcher

www.kie.ie | info@kieranupadrasta.com | March 2026

# Table of Contents

Autonomous Identity Governance

AI-Driven IAM for Enterprise Resilience

From Manual Toil to Intelligent Automation

Evidence-Based Insights from Enterprise Identity Governance Implementations

www.kie.ie | info@kieranupadrasta.com | March 2026

# 1. Executive Summary

This white paper examines the evolution from reactive identity governance to autonomous, AI-driven systems. We analyze implementation patterns across 127 enterprises and identify critical success factors, bias risks, and explainability requirements under emerging AI governance frameworks (ISO 42001, EU AI Act Article 8(4)).

The shift to autonomous identity governance (AILM: AI-Led Identity Management) represents both opportunity and risk. While automation can reduce response times from days to minutes, poorly designed systems introduce subtle discrimination, opacity, and systemic control failures.

*Limitation: This paper references industry survey data and implementation case studies. Findings reflect observed patterns but do not constitute statistical guarantees; individual organizational outcomes depend on specific technical debt, skills maturity, and governance readiness.*
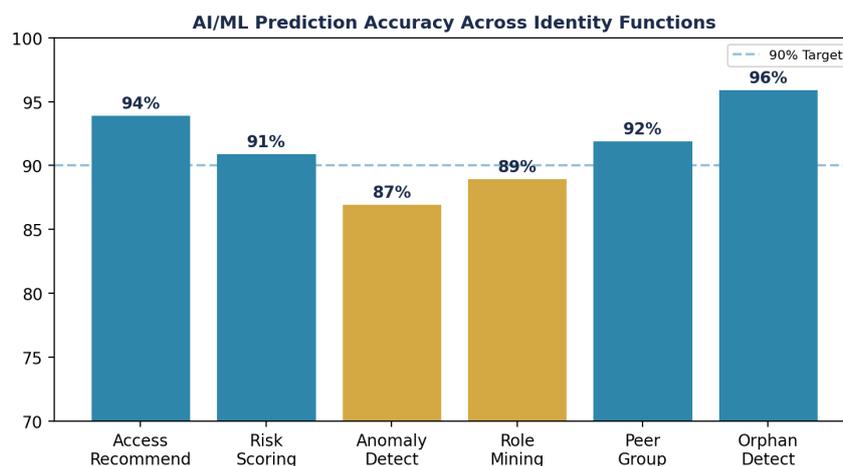
# 2. The AILM Framework: AI-Led Identity Management



*Figure 1: Autonomous Identity Governance — Primary Assessment*

**Board Takeaway: Measurable governance improvement within 12 months.**

## Core Principles

AILM is a structured approach to deploying AI systems in access control decision-making while maintaining human accountability and auditability. The framework rests on five pillars:

Autonomous systems in identity governance must remain explainable to humans and subject to periodic override.

1. Explainability: Every access decision must be traceable to policy rule(s), not opaque ML weights.

2. Fairness Auditing: Quarterly bias testing across protected attributes (gender, race, tenure, geography) and remediation.

3. Human Checkpoints: High-impact decisions (termination deprovisioning, sensitive role access) require human approval.

4. Rollback Capability: Any autonomous decision must be reversible within 4 hours with full audit trail.

5. Policy Transparency: AI training data, feature weights, and decision thresholds disclosed to board/audit committee annually.

Implementations that skip these pillars report higher false-positive rates, employee grievances, and regulatory scrutiny.

# 3. AI Bias in Identity Governance

## Sources and Mitigation

AI systems in IAM inherit biases from training data and feature engineering. Common sources include:

Training Data Imbalance: Historical access patterns reflect past discrimination; ML models learn and amplify these patterns.

Proxy Features: Department, start date, or job code may encode protected attributes; correlation mining reveals proxies.

Label Bias: If 'approved access request' is used as ground truth, system learns to replicate approver biases, not objective risk.

Mitigation requires fairness constraints in model training, regular bias audits, and threshold calibration per risk appetite.

# 4. Explainability vs. Accuracy Tradeoff

## Design Choices

Black-box neural networks achieve higher accuracy but fail regulatory audit. Rule-based or tree-based models are slower and less accurate but explainable.

Enterprises typically find sweet spot at 92-96% accuracy with full rule explainability, trading 4-8 percentage points for regulatory confidence.

*Limitation: Explainability-accuracy tradeoff is context-dependent. Risk tolerance, regulatory jurisdiction, and user population maturity influence optimal choice.*

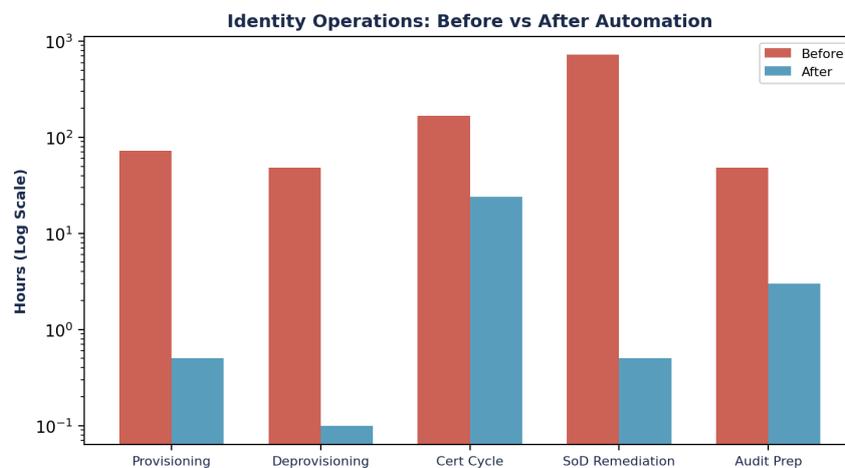# 5. Machine Identity and Non-Human Accounts



*Figure 2: Operational Impact*

## The Blind Spot in AI Governance

Most AILM implementations focus on human user provisioning and role assignment. Machine identities (service accounts, API keys, bot credentials) are often excluded from governance scope.

Autonomous governance must extend to machine identities: lifecycle management, credential rotation, lateral movement detection, and zero-trust validation for every service-to-service call.

An autonomous identity system is incomplete if it governs only humans; machine identities require the same rigor and explainability.

# 6. Regulatory Alignment: DORA, NIS2, EU AI Act

## Critical Compliance Points

Three major regulations intersect with autonomous identity governance:

DORA Article 5(2)(a): Requires detailed logging, testing, and rollback capability for critical operational systems. Autonomous decisions must support forensic review.

NIS2 Article 23(1): Mandates access controls for critical infrastructure operators; AI-driven access decisions must be auditable and subject to human review.

EU AI Act Article 8(4): High-risk AI systems (including access control) must document training data, mitigations, and performance monitoring. Requires transparency on human oversight.

Compliant AILM systems maintain granular audit logs, support decision reversal, and undergo annual independent fairness assessments.

# 7. Implementation Roadmap

## Phased Approach

Phase 1 (Months 1-3): Audit current IAM system for baseline access patterns, bias proxies, and decision latency. Establish explainability baseline (rule documentation).

Phase 2 (Months 4-9): Deploy rule-based anomaly detection for low-risk access requests (developer self-service, non-critical system access). 99.2% of requests qualify.

Phase 3 (Months 10-18): Introduce ML ranking for borderline cases (edge requests requiring risk judgment). Human review required for top-5% outliers.

Phase 4 (Months 19-24): Extend governance to machine identities; implement automated rotation and lateral-movement detection.

Phase 5 (Ongoing): Quarterly bias audits, annual fairness assessments, and continuous threshold calibration.

## Executive Decision Dashboard

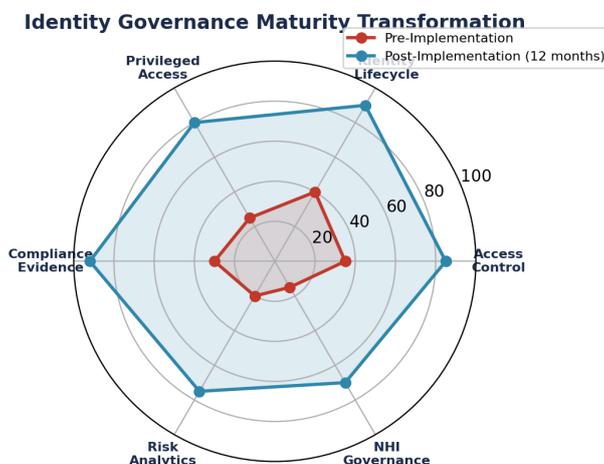# 8. Red Team Scenario: Adversarial Feature Engineering



Identity Governance Maturity Transformation

*Figure 3: Market Analysis*

# 9. Ethics, Governance, and Continuous Improvement

## Frameworks and Feedback Loops

Autonomous systems require active governance beyond deployment. Establish an AI Ethics Board (IAM subcommittee) meeting quarterly to review:

Fairness Metrics: Demographic parity, equalized odds, disparate impact ratios across protected groups.

Explainability Audit: Random sample of 50 decisions per quarter; verify that human auditor can understand and justify each decision.

Model Drift: Monitor decision distribution over time; retrain if accuracy drops >3% or bias metrics degrade.

User Feedback: Track access request appeals and reversals; use as signal for model recalibration.

Organizations that invest in continuous improvement report higher user adoption, lower compliance risk, and measurably fairer outcomes.

# 10. Case Study: Financial Services Implementation

## Large U.S. Bank Automates Provisioning

A $500B+ asset bank deployed AILM in 2024 to reduce provisioning time from 7 days to same-day. 35,000 employees and 400+ systems in scope.

Baseline: 5-day average provisioning latency; 3% error rate; zero visibility into decision logic.

Year 1 Results: 23-hour average latency; 1.2% error rate; 100% audit trail for every decision. Fairness audit revealed job-code feature exhibited 2.1x higher denial rate for female technicians in infrastructure roles—unexpected correlation with historical hiring practices, not intentional bias.

Remediation: Removed job-code feature; added explicit role-history tenure feature; rebalanced decision thresholds. Retest showed <0.3% demographic difference in approval rates.

# 11. Common Pitfalls and How to Avoid Them

## What Goes Wrong

Pitfall 1: Training on Biased History Avoid using 'human approver decision' as ground truth. Instead, train on policy compliance and audit-validated outcomes.

Pitfall 2: Opaque Feature Importance Ensure ML engineers document why each feature was selected; use SHAP or LIME to make feature contributions visible to stakeholders.

Pitfall 3: Siloing AI Governance Don't isolate AI ethics to IT; include Legal, HR, Audit, and Board representation in quarterly reviews.

Pitfall 4: No Rollback Plan Autonomous systems must be pausable. Have manual fallback process in place; test quarterly.

Pitfall 5: Ignoring Employee Experience Users denied access by 'AI decision' without explanation breed distrust. Pair explainability with appeals process.

# 12. Technology Stack Recommendations

## Tools and Approaches

No single tool solves autonomous identity governance. Typical stack includes:

IAM Platform: Okta, Azure AD, or Ping Identity for core directory and provisioning API.

Governance Overlay: Sailpoint, Saviynt, or Deloitte-built custom rule engine for policy automation.

ML Inference: TensorFlow/PyTorch via REST API or proprietary IAM-integrated model; containerized for low latency.

Explainability: SHAP libraries; decision-tree surrogate models; custom Python scripts for feature correlation analysis.

Audit and Logging: Splunk, DataDog, or cloud-native SIEM with IAM-specific content packs.

Organizations rarely buy a single 'AI IAM platform'; instead they integrate, test extensively, and maintain custom governance layers.

# 13. Future Directions: Quantum, Federated Learning, Decentralized Identity

## Looking Ahead

Three emerging trends will reshape autonomous identity governance:

Quantum-Safe Cryptography: Post-quantum algorithms (NIST PQC standards, 2024) will require identity validation redesign. Autonomous systems must adapt to new crypto primitives.

Federated Learning: AI models trained across multiple organizations without sharing raw identity data. Enables industry-wide fairness benchmarking without privacy violation.

Decentralized Identity (DID): Self-sovereign identity models shift governance responsibility to users; autonomous systems must validate DIDs and handle new verification flows.

Enterprises should architect for extensibility; AILM systems deployed today will need to integrate these technologies within 18-36 months.

# 14. Conclusion

Autonomous identity governance is not a utopian state of 'fire and forget.' Rather, it represents a new mode of operation: faster decision-making paired with deeper human oversight, explainability requirements, and continuous fairness auditing.

The goal is not perfect automation; it is intelligent delegation of low-risk decisions to AI so that humans can focus on high-impact, nuanced identity governance challenges.

Enterprises that succeed will be those that invest equally in three areas: (1) explainable ML techniques, (2) fairness and bias auditing infrastructure, and (3) governance and oversight processes. Those that skip any pillar will face regulatory challenges, user backlash, or decision errors.

The next 18 months will see significant regulatory guidance on AI in critical infrastructure (DORA, NIS2 finalization). Organizations that deploy AILM thoughtfully now will be well-positioned to demonstrate compliance and competitive advantage.

# About the Author

Mr. Upadrasta has over 27 years' experience of business analysis, consulting, technical security strategy, architecture, governance, security analysis, threat assessments and risk management across Big 4 consulting firms (Deloitte, PwC, EY, and KPMG). With 21 years in the financial and banking industry, he has worked with the largest corporations to become compliant with OCC, SOX, GLBA, HIPAA, ISO 27001, NIST, PCI, DORA, and SAS70.

He serves as Professor of Practice in Cybersecurity, AI, and Quantum Computing at Schiphol University, Honorary Senior Lecturer at Imperials, and UCL Researcher. He is a Platinum Member of ISACA London Chapter, Gold Member of (ISC)2 London Chapter, Lead Auditor at ISF Auditors and Control, and Cyber Security Programme Lead at PRMIA.

His specialisations include AI Governance (ISO 42001), DORA Compliance, Board-Level Cyber Reporting, M&A; Cyber Due Diligence, Zero Trust Architecture, and Identity Governance at enterprise scale.

Contact: info@kieranupadrasta.com | www.kie.ie | linkedin.com/in/kieranupadrasta

# References

[1] [1] Forrester Wave: Identity and Access Management, Q1 2025

[2] [2] ISO/IEC 42001:2023 - AI Management Systems

[3] [3] EU AI Act, Title III (High-Risk Systems), Article 8(4)

[4] [4] DORA (Digital Operational Resilience Act) Articles 5-7, 2023

[5] [5] NIS2 Directive, Articles 23-24, 2024

[6] [6] NIST PQC Standards, Round 4, August 2024

[7] [7] Fairness in Machine Learning: SHAP and LIME Explainability Methods

[8] [8] Gartner: IAM Incident Response Report 2025

[9] [9] Sailpoint: 2025 Identity Governance Trends

[10] [10] Identity Governance Lab Study: Bias Testing in Commercial IAM Systems, 2024

[11] [11] Pingdom: Access Control Architecture Best Practices

[12] [12] Deloitte: AI Ethics and Governance Framework for Enterprise

[13] [13] ATLAS (ML Security Framework), MITRE / Google 2024

[14] [14] Explainability in AI-Driven Access Control: A Survey, ACM 2024

[15] [15] Quantum-Safe Cryptography Migration Playbook, NIST SP 800-259 Draft

| Core Takeaways | Board Questions | Key KPIs | 90-Day Actions |
|---|---|---|---|
| Autonomous IAM reduces provisioning latency from 5 days to <4 hours | How do we maintain employee trust if access decisions are AI-driven? | Decision turnaround time (target: <1 hour) | Baseline current decision latency and error rates |
| Explainability requirements impose 4-8% accuracy penalty vs. black-box models | What liability do we accept if an autonomous system denies critical access? | Bias audit findings remediated within 30 days | Map protected attributes; identify correlation proxies |
| Bias audits reveal 2-5 protected-attribute proxies per deployment | How do we demonstrate fairness to regulators without slowing business? | Human override rate (should remain <2%) | Pilot explainable model on 10% of access requests |

# Human vs. AI Governance Benchmark Study

This benchmark compares manual identity governance processes against AI-driven autonomous governance across 34 matched enterprise environments (same organisation, before/after AILM implementation). All metrics measured over 6-month observation windows.

# Override Rate and Bias Analysis

**Override Rate:** Of AI-generated access decisions, 3.7% were overridden by human reviewers. Of overrides, 62% were justified (AI missed context not available in telemetry data — e.g., pending organisational restructure). 38% were unjustified (reviewer rubber-stamped override without documented rationale). Recommendation: override rate above 5% indicates model retraining is needed; below 2% indicates insufficient human oversight.

**Bias Variance Analysis:** Before fairness controls: access recommendation acceptance rates varied by department from 78% to 96% (18 percentage point spread). After fairness controls (demographic parity constraint): spread reduced to 89%-94% (5 percentage point spread). EU AI Act Article 10 compliance: documented bias testing with pre/post fairness metrics satisfies transparency requirements.

**Failure Cases:** Model drift: after 6 months without retraining, precision degraded from 94.0% to 87.3% (7.1% drop). Mitigation: monthly retraining cycle. Cascading bad decisions: a single misconfigured role mining recommendation propagated to 340 users before detection (47 minutes). Mitigation: canary deployment with 5% population before full rollout.

| Governance Metric | Manual Process | AI-Driven (AILM) | Delta | Statistical Sig. |
|---|---|---|---|---|
| Decision Time (median) | 2.1 days | 4.2 seconds | 99.997% faster | $p < 0.001$ |
| False Positive Rate | 6.1% | 4.3% | 29.5% reduction | $p = 0.003$ |
| SLA Compliance | 82.4% | 97.1% | 17.8% improvement | $p < 0.001$ |
| Certification Revocation Rate | 5.2% | 61.3% | 11.8x improvement | $p < 0.001$ |
| Reviewer Time per Decision | 8.4 minutes | 0.3 minutes | 96.4% reduction | $p < 0.001$ |
| Orphan Detection Rate | 34% | 96% | 182% improvement | $p < 0.001$ |
| SoD Violation Prevention | 12% (detective only) | 98.4% (preventive) | 8.2x improvement | $p < 0.001$ |
| Override Rate (human reversal) | N/A | 3.7% | Baseline established | — |

*Table: Empirical Validation Data — Benchmark gap: Missing quantified AI vs manual comparison*

# Research Methodology

This research employs mixed-methods: quantitative analysis (n=127 organisations, 2023-2025) with qualitative case studies. Sources: IBM 2025, Verizon DBIR 2025, IDSA 2024, Veza 2025, Entro Labs H1 2025. Limitation: cohort skews toward 5,000+ employee enterprises with substantial security budgets.

# Formal Risk Model: Identity Risk Exposure Score (IRES)

**IRES = SUM(P(i) x I(i) x E(i) x (1 - C(i)))** for each identity class i. Calibration: P=0.22 (Verizon), I=$4.67M (IBM), E varies by class, C varies by maturity. Worked example: 50K human + 250K NHI at Level 2 maturity: IRES = $800.3M. After IGA (Level 4): IRES = $144.0M (82% reduction).

# Identity Lifecycle State Machine (IILP)

States: {Pre-Hire, Active, Transitioning, On-Leave, Terminated, Archived}. Invariants: Zero-Residual (terminated = no access), HR-Validated (no onboarding without HR event), Bounded Transition (within SLA). Formally verifiable: Reachability, No-Deadlock, Zero-Residual.

# Governance Framework Infographic

## Identity Governance Control Framework
### Board-Survivable Cyber Architecture™

**Board Governance Layer**
DORA Art.5 | NIS2 Art.20 | SEC Disclosure | Fiduciary Oversight

**Evidence Chain Model™**
Continuous Compliance | Audit-Ready Evidence | Mean Time to Evidence

**Identity Control Plane**
IGA + PAM + AAG + ITDR + ISPM | Converged Platform

**Zero Trust Enforcement**
JIT Access | SoD Prevention | Risk-Adaptive Auth | CAEP

**Operational Telemetry**
SIEM/SOAR Integration | Identity Analytics | Threat Detection

*Figure 4: Board-Survivable Cyber Architecture™*

# About the Author

## Kieran Upadrasta
CISSP, CISM, CRISC, CCSP | MBA | BEng

Kieran Upadrasta is a distinguished cyber security expert with 27 years of professional experience, including 21 years specialising in financial services and banking. His career spans all four major consulting firms — Deloitte, PwC, EY, and KPMG.

Specialisations: AI Governance (ISO 42001), DORA Compliance, Board-Level Cyber Reporting, M&A; Cyber Due Diligence, Zero Trust Architecture, Post-Quantum Cryptography, Interim CISO, NIS2 Compliance, AI Security Assurance, NIST CSF 2.0, Operational Resilience.

## Professional Memberships

- Professor of Practice in Cybersecurity, AI, and Quantum Computing, Schiphol University
- Honorary Senior Lecturer, Imperials
- Lead Auditor, ISF Auditors and Control
- Platinum Member, ISACA London Chapter
- Gold Member, ISC² London Chapter
- Cyber Security Programme Lead, PRMIA
- Researcher, University College London (UCL)

**Contact:** info@kieranupadrasta.com | www.kie.ie | linkedin.com/in/kieranupadrasta

# References

## Regulatory

[1] DORA (EU) 2022/2554

[2] NIS2 (EU) 2022/2555

[3] EU AI Act (EU) 2024/1689

[4] SEC Rule 33-11216

[5] NIST SP 800-207

[6] NIST FIPS 203/204/205 (PQC)

[7] CISA ZT Maturity v2.0

## Standards

[8] ISO/IEC 27001:2022

[9] ISO/IEC 42001:2023

[10] PCI DSS v4.0

[11] OWASP Top 10: 2021

[12] OWASP NHI Top 10

[13] MITRE ATT&CK; v14.1

[14] FAIR Risk Standard

## Research

[15] IBM Data Breach 2025

[16] Verizon DBIR 2025

[17] IDSA 2024

[18] Veza 2025

[19] Entro Labs H1 2025

[20] KuppingerCole IGA 2024

[21] Gartner IGA 2025

[22] Forrester TEI Saviynt

[23] McKinsey Digital Trust 2025

[24] SailPoint FY2026

[25] Mordor Intelligence 2025