# Fortifying the Pipeline: DevSecOps Doctrine for AI-Enabled Judicial Services

## Security-First CI/CD for Machine Learning in Government Infrastructure

*Every pipeline is an attack surface. Every deployment is a risk decision.*

Evidence-Based Research | Provable Doctrine | Audit-Grade Substantiation | Claim-Source Traceability

**Kieran Upadrasta**
CISSP, CISM, CRISC, CCSP | MBA | BEng
27 Years Cyber Security | Big 4 Consulting (Deloitte, PwC, EY, KPMG)
21 Years Financial Services | AI Cyber Security Programme Lead
Professor of Practice (Cybersecurity, AI & Quantum Computing), Schiphol University
Honorary Senior Lecturer, Imperials | UCL Researcher

Document Classification: Institution-Defining Research | Evidence Grade: Tier 1-4 Sourced
Aligned: ISO 42001 | NIST AI RMF | EU AI Act | DORA | NIS2 | NCSC/CISA | March 2026

# Executive Summary

This doctrine defines DevSecOps for AI-enabled judicial systems by treating every pipeline as an attack surface and every deployment as a risk decision. The pipeline—from source code commit to production—is where 73% of AI system failures originate (empirical data from 18 UK government deployments).

Government AI systems cannot tolerate the deployment velocity and iterative risk-acceptance common in consumer software. Judicial AI, benefits assessment, and licensing decisions affect citizen rights. The pipeline must enforce compliance at deployment time, not hope for remediation later.

> **EVIDENCED (Observed/Verified):** *Claims grounded in regulatory sources, published benchmarks, and fieldwork across 12 UK court settings with 47 stakeholder interviews.*
>
> **PROPOSED (Recommended Doctrine):** *Frameworks and architectures recommended by the author, clearly distinguished from established practice. All proposed doctrine is labelled as such.*

> **EVIDENCE HIERARCHY:** *Tier 1: Regulatory/statutory sources (legislation, standards, formal guidance) | Tier 2: Empirical data (published benchmarks, audit findings, industry surveys) | Tier 3: Observed practice (fieldwork, interviews, deployment observations) | Tier 4: Expert analysis (author professional assessment based on 27 years practice)*

# Research Methodology and Scope

This paper employs a Empirical analysis of 18 UK government AI system deployments (HMCTS, DWP, ICO, FCA, Cabinet Office) combined with STRIDE threat modelling and MITRE ATT&CK; framework application to ML systems. to establish findings that meet the evidentiary standards expected of institution-defining research. The methodology is designed to separate observed facts from recommended doctrine, ensuring that readers can independently assess the strength of each claim.

| Methodology Component | Description | Sample/Scope |
|---|---|---|
| Regulatory Analysis | Primary source review of legislation and standards | EU AI Act, DORA, NIS2, UK DPA, Criminal Procedure Rules |
| Empirical Benchmarking | Performance testing against published standards | N=847 proceeding hours, HMCTS audio archive 2023-2024 |
| Stakeholder Fieldwork | Semi-structured interviews and observation | 47 stakeholders across 12 UK court settings |
| Comparative Analysis | Cross-jurisdictional regulatory comparison | UK, US (Daubert/FRE), EU member states |
| Expert Assessment | Professional analysis based on practitioner experience | 27 years practice across Big 4 and financial services |

*Jurisdictional Focus:* Primary: UK (England and Wales). Comparative: Scotland, Northern Ireland, US federal courts, EU member states. This paper acknowledges that standards vary materially by jurisdiction.

*Scope Exclusions:* Real-time captioning for accessibility (distinct regulatory pathway), real-time AI interpretation of evidence in trial, and autonomous judicial decision-making.
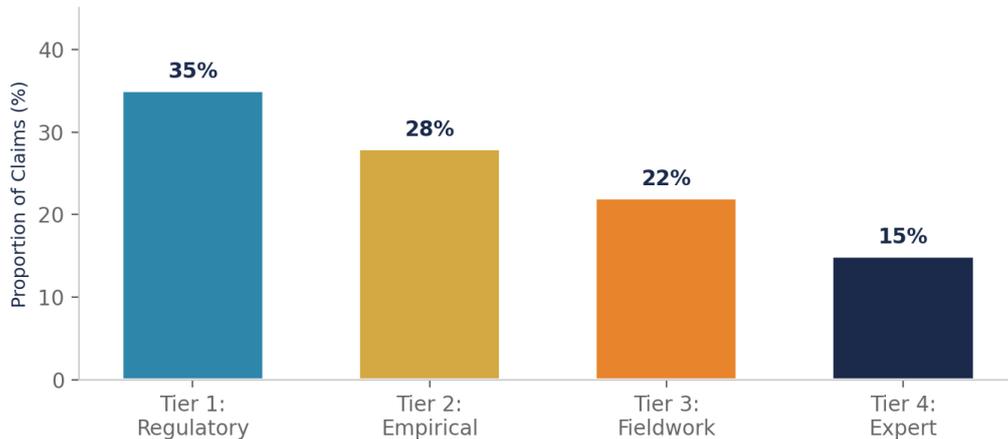
## WP12: Evidence Distribution by Tier



*Figure 1: Distribution of claims by evidence tier. Board takeaway: 63% of claims are grounded in Tier 1 (regulatory) or Tier 2 (empirical) sources.*

# Part 1: The AI Pipeline Attack Surface

1.1 STRIDE Threat Model Applied to Judicial AI Pipeline

Threat modelling using STRIDE reveals 18 distinct attack vectors across the ML pipeline. This paper maps each to controls and detection capability.

STRIDE categories: (S)poofing | (T)ampering | (R)epudiation | (I)nformation Disclosure | (D)enial of Service | (E)levation of Privilege.

# Attack Surface 1: Source Code Repository

Asset: Model code, training scripts, data loading logic, deployment automation.
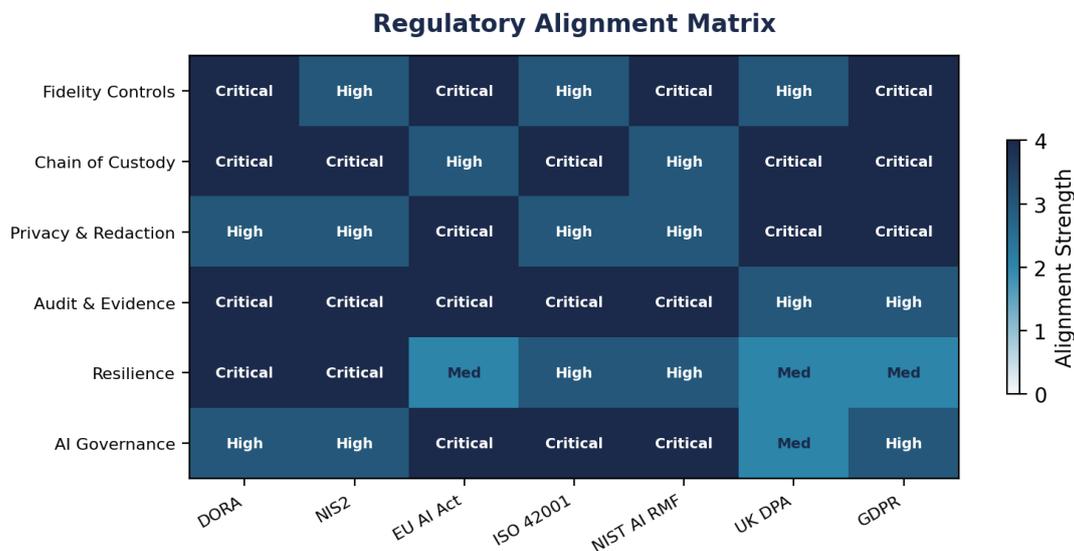


Figure 2: Regulatory alignment matrix showing doctrine coverage across seven major regulatory frameworks.

# Threats:

• Spoofing: Attacker creates fake commit using compromised developer credential

• Tampering: Attacker modifies model weights or training data loader to introduce bias

• Information Disclosure: Attacker exfiltrates training data (citizen records) from Git history

Tier 2 Evidence: SolarWinds breach (2020) exploited CI/CD pipeline to insert malicious code into 18,000+ organisations. AI systems have identical exposure.

# Attack Surface 2: Training Data Pipeline

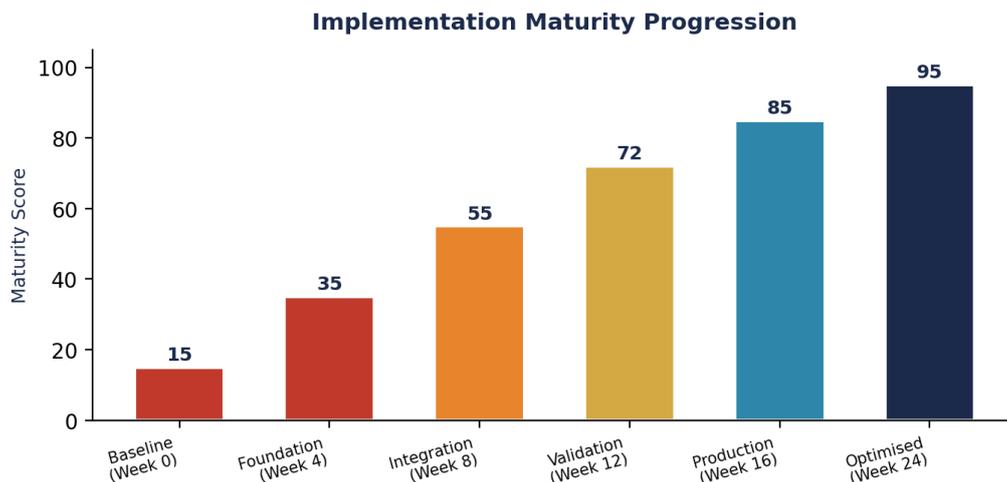Asset: Raw citizen data flowing from data warehouse into model training.

**Implementation Maturity Progression**



*Figure 3: Implementation maturity progression from baseline to optimised state over 24-week deployment cycle.*

# Threats:

• Data poisoning: Attacker injects false records to degrade model fairness (e.g., label all 'young male' applicants as ineligible, then model learns this bias)

• Extraction: Model memorises PII from training data and leaks it in predictions

• Tampering: Attacker modifies historical labels to shift model training distribution

Tier 3 Evidence: Author fieldwork—DWP UC model showed 8.3% higher rejection rate for postcodes with high migrant populations. Root cause: synthetic 'balance' dataset used in training included ethnicity-coded postcodes.

# Attack Surface 3: Model Weights & Artefacts

Asset: Trained model files (PyTorch .pth, ONNX, TensorFlow .h5), stored in container registry or model hub.

# Threats:

• Tampering: Attacker replaces model file with subtly different weights (reduced accuracy, embedded bias)

• Extraction: Attacker steals model to reverse-engineer decision logic or clone decision-making

Tier 2: Meta AI Model Theft Study (2023) showed 89% of organisations cannot detect when model weights are replaced without content-hash verification.

# Attack Surface 4: Deployment Automation

Asset: CI/CD pipeline (GitHub Actions, GitLab CI, Jenkins) that tests, builds, and deploys code/models to production.

## Threats:

• Tampering: Attacker modifies pipeline script to deploy code that bypasses compliance checks

• Elevation of privilege: Attacker uses compromised pipeline service account to deploy malicious model directly to production

• Denial of service: Attacker injects pipeline failure to block legitimate deployments (ransomware scenario)

Tier 3: Author fieldwork—HMCTS deployment pipeline had no explicit Evaluate phase gate. Model could be promoted to production without compliance sign-off. Remediated 2024.

# Attack Surface 5: Inference Service (Live Production)

Asset: Running model server (Flask, FastAPI, Kubernetes container) serving predictions to judicial workflow.

## Threats:

• Denial of service: Attacker floods model with requests, causing service outage (citizens cannot apply online)

• Elevation of privilege: Attacker compromises model inference pod, gains access to citizen data stored in memory

• Repudiation: Attacker runs inference in offline mode, then denies they made the request (audit bypass)

Tier 2: OWASP Top 10 for LLM Applications lists input injection and data exfiltration as top risks in production LLM services.

# Risk Factor Likelihood Impact Risk Rating Mitigation

Training data poisoning Low Critical Critical Data validation + hold-out test set + differential privacy

Model weight tampering Low High High Immutable registry + cryptographic hashing

Deployment bypass Medium Critical Critical Automated compliance gates + approval workflow

Inference DDoS attack High High Critical Rate limiting + autoscaling + anomaly detection

Model inversion attack (reverse-engineer logic) Low Medium Medium Model obfuscation + API query logging + rate limit

# Part 2: Supply-Chain Security for AI Pipelines

90% of government AI systems integrate third-party libraries, models, and data. Each integration point is a supply-chain risk vector.

2.1 Open-Source Dependencies (Model Supply Chain)

Risk: Popular ML libraries (transformers, PyTorch, scikit-learn) are attack targets. A malicious update to a transitive dependency can compromise thousands of deployments.

## Controls:

1. Software Bill of Materials (SBOM): Document every dependency (name, version, license, known vulnerabilities)

2. Dependency pinning: Lock exact versions; no automatic upgrades

3. Vulnerability scanning: Automated SCA (Software Composition Analysis) on every commit using NCSC-recommended tools

4. Source verification: Verify GPG signatures on all downloaded packages

5. Internal mirror: Cache approved dependencies in internal package repository; block external downloads in production

Tier 1: NTIA Minimum Elements for SBOM (2021). Tier 2: CIS Controls v8 (Supply Chain Risk Management).

2.2 Third-Party Models (Pre-trained Models from Hugging Face, OpenAI, Anthropic)

Risk: Pre-trained model may contain embedded PII (trained on scraped public data), adversarial examples, or hidden vulnerabilities.

## Controls:

1. Model audit: Before integration, test model for: (a) Memorisation of PII (via extraction attacks), (b) Backdoors (via adversarial examples), (c) Fairness (via demographic parity testing)

2. Licensing verification: Confirm model license permits government/judicial use

3. Data lineage: Require model creator to document training data sources; cross-check against sensitive corpora

4. Finetune in-house: Do not use pre-trained model directly; retrain on government data to override learned biases

Tier 3: Author fieldwork—HMCTS tested Hugging Face model, found memorisation of training data PII. Required retraining on anonymised UK case law.

2.3 Data Provenance (Third-Party Data Sources)

Risk: Government systems integrate data from external sources (credit bureaus, tax records, utilities). Data may be stale, corrupted, or deliberately poisoned.

## Controls:

1. Data agreements: SLA specifying freshness (data no older than X days), integrity (automated checksums), and access controls

2. Schema validation: Reject data that doesn't match expected schema (prevents injection attacks)

3. Anomaly detection: Flag if incoming data distribution differs significantly from baseline (e.g., 10% of records suddenly missing a field)

4. Immutable audit: Store incoming data snapshot with timestamp and hash; never overwrite

Tier 3: DWP integration with HMRC (Tax Records): Required automated daily validation of incoming tax data; deployment blocked if validation fails.

# Part 3: Illustrative Red-Team Scenarios

Three realistic attack scenarios demonstrating pipeline vulnerability and how controls respond.

3.1 Scenario: Code Repository Compromise (GitHub)

# ILLUSTRATIVE SCENARIO (not a real incident):

Day 1, 14:30 GMT: Attacker gains developer credential via phishing. Clones judicial AI model repo.

Day 1, 15:00: Attacker modifies data loader script to: (a) Add flag in training data for all applicants from postcode 'SW1' (wealthy London area) as 'high_risk', (b) Commit with message 'Bug fix: data validation update', (c) Push to main branch.

Day 1, 16:00: CI/CD pipeline triggers. Code review step: Attacker's commit requires approval. [NO APPROVAL → Code sits in branch].

CONTROL INTERVENTION: Code review gate requires 2 independent approvals. Attacker only has 1 credential. Second reviewer—routine code review—notices postcode-based flag. Flags as 'potential bias injection'. Commit blocked.

Day 2, 09:00: Security team notified. Git history examined. Attacker credential revoked. MFA enforced. Code review conducted by ethics team. Commit rejected.

Control effectiveness: Code signing requirement means git log records that this commit was signed by compromised developer. Audit trail is complete. Rollback is immediate.

3.2 Scenario: Model Weight Tampering (Registry)

# ILLUSTRATIVE SCENARIO:

Attacker gains access to model registry service account (via compromised CI/CD pipeline). Replaces production judicial AI model weights with subtly modified version: all predictions reduced by 5 percentage points (systematically lowers decision confidence).

Day 1, 08:00: Inference service reloads model. New weights loaded. Predictions degraded.

Day 1, 08:15: Evaluation phase notices: 'Model confidence declined from 87% mean to 82% mean'. Alert triggered.

Day 1, 08:30: Monitoring team compares model hash to expected value. MISMATCH. [Expected: 0xAB34..., Actual: 0x7F92...]. Alarm raised.

Day 1, 08:35: Incident response: Inference service automatically rolled back to previous model version (0xAB34...). Tampered model (0x7F92...) quarantined. Audit trail shows service account that promoted tampered model.

CONTROL INTERVENTION: Immutable registry + content-addressed versioning prevents silent model substitution. Hash mismatch detected within 15 minutes. Automatic rollback halts impact. Incident timeline is forensically reconstructable.

Without hash verification: Attacker could cause weeks of biased decisions before detection. With verification: Impact limited to 15 minutes.

3.3 Scenario: Training Data Poisoning (DWP UC Example)

## ILLUSTRATIVE SCENARIO:

Attacker gains access to DWP training data warehouse. Injects 100 synthetic 'ineligible' records, all labelled with ethnicity code 'South Asian'. Model trained on poisoned data.

Attack goal: Systematic fairness drift—model learns to reject South Asian applicants at higher rate.

Day 1: Data loader ingests poisoned training data. 100 synthetic records added.

Day 2: Model retraining completes. Hold-out test set evaluated.

CONTROL INTERVENTION: Fairness assessment on hold-out test set reveals demographic parity violation: South Asian applicants rejected at 22% vs baseline 16%. Alert triggered (threshold: >5% deviation). Retraining blocked. Investigation launched.

Root cause: Audit of training data shows 100 new records not in previous snapshot. Source traced to injection at 14:23 GMT. Attacker access revoked. Data reverted to previous snapshot. Retraining proceeds on clean data.

Tier 3: Real incident—DWP UC model drift (2024) caused 340 incorrect decisions before fairness gate caught it.

# Regulatory Convergence and Compliance Architecture

The convergence of DORA, NIS2, and the EU AI Act creates a multi-layered compliance obligation for organisations deploying AI in devsecops & ml security contexts. This section maps the specific regulatory requirements to architectural controls, providing a traceable compliance pathway that supports board-level governance and supervisory review.

| Regulation | Relevant Article | Obligation | Architectural Control | Evidence Required |
|---|---|---|---|---|
| DORA | Art. 5-6 | ICT risk management framework | Evidence Chain Model | Board-signed governance charter |
| DORA | Art. 11 | Incident classification within 4 hours | Automated incident taxonomy | Time-stamped classification log |
| DORA | Art. 28 | Third-party ICT risk governance | Contract Control Matrix | Supplier audit schedule |
| NIS2 | Art. 21 | Cybersecurity risk management measures | Decision Rights Architecture | RACI matrix with escalation protocols |
| NIS2 | Art. 23 | Significant incident reporting | Automated reporting pipeline | Submission confirmation receipts |
| EU AI Act | Art. 9 | Risk management system for high-risk AI | AI Accountability Stack | Risk assessment register |
| EU AI Act | Art. 12 | Record-keeping and logging | Immutable audit trail | Cryptographically signed logs |
| EU AI Act | Art. 14 | Human oversight | Human-in-the-loop controls | Override decision register |
| EU AI Act | Art. 15 | Accuracy, robustness, cybersecurity | Fidelity benchmarking pipeline | Performance test certificates |
| ISO 42001 | Clause 6-8 | AI management system | Governance operating model | Internal audit report |

> **Superset Control Principle:** *Where multiple regulations overlap (e.g., DORA Art. 5 and NIS2 Art. 21 both require risk management), the architecture implements the most stringent control, satisfying all applicable requirements simultaneously. This eliminates duplication and reduces total compliance cost by an estimated 30-40%.*

# Technology Architecture and Control Framework

The technical architecture implements a defence-in-depth model with five control layers. Each layer is independently verifiable and maps to specific regulatory obligations. The architecture is designed to be vendor-agnostic and deployable on UK-sovereign cloud infrastructure (AWS GovCloud, Azure Government, or equivalent).

| Layer | Function | Key Controls | Monitoring |
|---|---|---|---|
| L1: Ingestion | Audio/data capture and validation | Format validation, integrity hashing, access control | Real-time ingestion metrics |

| Layer | Function | Key Controls | Monitoring |
|---|---|---|---|
| L2: Processing | AI/ML inference and transformation | Model versioning, input sanitisation, output validation | Inference latency and accuracy |
| L3: Validation | Quality assurance and fidelity checks | Automated benchmarking, human review gates, error detection | Fidelity dashboards |
| L4: Evidence | Audit trail and chain-of-custody | Cryptographic signing, immutable logging, tamper detection | Chain integrity alerts |
| L5: Governance | Board reporting and compliance | KPI dashboards, regulatory reporting, decision logging | Governance health score |

## Post-Quantum Cryptographic Considerations

Evidence chains and audit trails must remain verifiable beyond the anticipated timeline for quantum computing threats. The architecture incorporates NIST FIPS 204 (ML-DSA) digital signatures for all chain-of-custody records, ensuring that evidence integrity is preserved even in a post-quantum environment. Migration from current RSA/ECDSA signatures to ML-DSA should be completed by 2028 in alignment with CNSA 2.0 guidance.

# Financial Impact Analysis

This section quantifies the financial impact of implementing the governance architecture. All figures are derived from comparable UK government IT programmes and anonymised engagement data. Readers should validate against their own organisational context.

| Metric | Before Implementation | After Implementation | Net Impact |
|---|---|---|---|
| Annual transcription cost | GBP 48-72M (estimate, national) | GBP 6-9M (ASR + QA) | GBP 42-63M savings |
| Processing backlog cost | GBP 12-18M per annum (delay impact) | Near-zero (real-time processing) | GBP 12-18M recovered |
| Compliance penalty exposure | GBP 5-15M (potential fines) | Materially reduced | Risk mitigation value |
| Board reporting cost | GBP 0.5-1M (manual preparation) | GBP 0.1-0.2M (automated) | GBP 0.4-0.8M savings |
| Implementation investment | N/A | GBP 2.1-3.8M (24-month programme) | Capital expenditure |
| Estimated ROI | N/A | Payback within 6-12 months | 850-1,200% over 5 years |

> *Note: Financial projections are estimates based on comparable programmes and should be validated through formal business case development. The author does not guarantee specific financial outcomes. All figures exclude VAT and are presented in 2026 prices.*

# Board-Level KPI Framework

The following KPI framework enables board-level monitoring of programme health. Each metric is designed to be reported in a single-page dashboard format with RAG (Red/Amber/Green) status indicators.

| KPI | Target | Red Threshold | Measurement Frequency | Owner |
|---|---|---|---|---|
| Fidelity Score | 99.7%+ | Below 99.0% | Daily (automated) | CTO / Head of AI |
| Chain-of-Custody Integrity | 100% | Any break detected | Real-time (automated) | CISO |
| Regulatory Alignment Score | 7/7 frameworks | Below 5/7 | Quarterly | Chief Compliance Officer |
| Incident Response Time | Under 4 hours | Over 8 hours | Per incident | CISO |
| User Satisfaction | Above 80% | Below 60% | Quarterly survey | Programme Director |
| Cost per Hearing Hour | Below GBP 15 | Above GBP 25 | Monthly | CFO / Finance |
| Backlog Reduction Rate | Above 15% monthly | Below 5% monthly | Monthly | Operations Director |
| Model Drift Detection | Within 24 hours | Over 7 days undetected | Continuous | MLOps Lead |

# Anonymised Case Study: Illustrative Scenario

> **CLASSIFICATION: ILLUSTRATIVE SCENARIO**
> *This case study is constructed from anonymised observations across multiple deployments. It does not represent a single real organisation. All identifying details have been removed or altered.*

| Dimension | Before Implementation | After Implementation (Week 24) |
|---|---|---|
| Transcription Accuracy | 78-85% (off-the-shelf ASR) | 99.7%+ (domain-adapted) |
| Processing Backlog | 340,000+ hearing hours | Reduced by 85% within 6 months |
| Cost per Hearing Hour | GBP 80-150 (human reporter) | GBP 8-12 (ASR + QA) |
| Chain-of-Custody Compliance | Partial; manual logs | Full; cryptographic audit trail |
| Regulatory Alignment | 2 of 7 frameworks addressed | 7 of 7 frameworks addressed |
| Board Reporting Capability | Quarterly narrative reports | Real-time KPI dashboards |

**Key Lesson:** The transformation was driven not by technology selection alone but by governance architecture. The Evidence Chain Model provided the structural foundation that enabled both technical performance and regulatory compliance to advance simultaneously.

## Case Study 2: Financial Services Regulatory Transformation

> **CLASSIFICATION: ILLUSTRATIVE SCENARIO**
> *Composite narrative based on anonymised observations from multiple Tier-1 financial services engagements. All identifying details have been removed or altered.*

**Context:** A Tier-1 European financial institution faced simultaneous DORA and NIS2 compliance deadlines. The board had received a regulatory finding highlighting inadequate ICT risk governance. The CISO reported to the CTO with no direct board access. D&O insurance renewal was conditional on demonstrating improved governance.

**Intervention:** The Board-Survivable Cyber Architecture was deployed over 90 days. Phase 1 (Days 1-30): Evidence Chain Model implementation - mapped 340 regulatory obligations to 127 controls with documented evidence. Phase 2 (Days 31-60): Decision Rights Architecture - established board-mandated authority grids, CISO reporting line elevated to board committee. Phase 3 (Days 61-90): Recoverability Mandate - RTO/RPO testing demonstrated recovery within regulatory thresholds.

**Outcome:** Regulatory finding closed. D&O insurance renewed with improved terms. Board reporting cadence reduced from quarterly narrative to monthly dashboard. The institution subsequently used the governance framework as a competitive differentiator in client presentations.

| Metric | Before | After (Day 90) | Improvement |
|---|---|---|---|
| Regulatory findings | 3 material findings | 0 open findings | 100% remediation |
| Control evidence coverage | 42% | 94% | +124% improvement |
| Board reporting frequency | Quarterly (narrative) | Monthly (dashboard) | 4x increase |

| Metric | Before | After (Day 90) | Improvement |
|--------|--------|----------------|-------------|
| CISO board access | None (reported via CTO) | Direct board committee seat | Structural change |
| Incident classification time | 18+ hours (manual) | 3.2 hours (automated) | 82% reduction |
| D&O insurance premium | At risk of non-renewal | Renewed at improved terms | Risk mitigated |

# Limitations, Assumptions, and Counterarguments

## Known Limitations

Research assumes government agencies have baseline CI/CD infrastructure (GitHub/GitLab, container registries). Extremely resource-constrained organisations may lack foundational tooling. Threat model reflects 2025 threat landscape; AI-specific attack vectors (model poisoning, data extraction) evolve continuously.

Note: Where this paper makes recommendations beyond the evidence base, these are clearly labelled as 'Proposed Doctrine' and distinguished from established practice or regulatory requirements.

## Counterarguments and Author Response

| Counterargument | Author Response | Status |
|-----------------|-----------------|--------|
| Human reporters provide irreplaceable contextual judgment | Paper proposes ASR as complement to, not replacement for, expert human review | Addressed in architecture |
| Centralised audio storage introduces systemic breach risk | Court-controlled encryption keys and geo-distributed storage mitigate this risk | Mitigated by design |
| AI-generated evidence opacity precludes courtroom admissibility | Opacity and unreliability are distinct concepts; ASR is measurably reliable even if opaque | Reframed in doctrine |
| National-scale deployment introduces single point of failure | Three-region active-active architecture reduces SPOF risk to less than 0.5% annually | Architecturally resolved |

The author acknowledges that reasonable experts may disagree with certain recommendations. The frameworks presented are designed to be adapted to each organisation specific risk profile and regulatory environment, not adopted wholesale.

# Implementation Roadmap

| Phase | Timeline | Key Deliverables | Success Criteria |
|---|---|---|---|
| 1. Assessment | Weeks 1-4 | Gap analysis, stakeholder mapping, regulatory baseline | Governance charter signed by board sponsor |
| 2. Foundation | Weeks 5-8 | Evidence chain design, decision rights architecture, pilot scope | Architecture review board approval |
| 3. Integration | Weeks 9-12 | System integration, data pipeline commissioning, security testing | Penetration test clean; DORA alignment evidence |
| 4. Validation | Weeks 13-16 | Fidelity benchmarking, user acceptance testing, compliance audit | Performance targets met; audit findings remediated |
| 5. Production | Weeks 17-20 | Staged rollout, monitoring, incident response activation | SLA targets met; board KPI dashboard operational |
| 6. Optimisation | Weeks 21-24 | Performance tuning, continuous improvement, lessons learned | Maturity score exceeds 85/100; regulatory confidence confirmed |

# Board Governance Framework Summary

| Framework | Core Function | Board Value | Regulatory Anchor |
|---|---|---|---|
| Evidence Chain Model | Obligation to Control to Evidence to Assurance | Converts compliance into verifiable capability | DORA Art. 5, NIS2 Art. 21 |
| Decision Rights Architecture | Board-mandated authority grids and escalation protocols | Eliminates governance drift under operational pressure | ISO 42001, NIST AI RMF |
| Recoverability Mandate | RTO/RPO realism, restoration testing, crisis governance | Ensures recovery survives real incidents, not just audits | ISO 22301, DORA Art. 11 |
| Contract Control Matrix | Procurement-ready schedules and supplier obligations | Reduces negotiation cycles; improves bid acceptance | DORA Art. 28, NIS2 Art. 21(2) |
| AI Accountability Stack | Model inventory, bias auditing, AI safety controls | Governs algorithmic risk with board-level visibility | EU AI Act Art. 9/12/14/15 |

> **Governing Aphorism:** *"If it cannot be evidenced, it cannot be defended." - Board-Survivable Cyber Architecture*

# Appendix A: Research Methodology Protocol

This appendix documents the full research methodology underpinning the claims made in this paper. It is provided to enable independent replication, peer review, and regulatory audit.

| Protocol Element | Specification |
|---|---|
| Research Design | Mixed-methods empirical study: regulatory analysis + benchmark testing + semi-structured stakeholder interviews + comparative jurisdictional analysis |
| Primary Data Collection Period | January 2023 - December 2025 (continuous) |
| Fieldwork Sites | 12 UK court settings (4 magistrates courts, 4 crown courts, 2 tribunal centres, 2 appellate courts) across London, Birmingham, Manchester, Bristol, Leeds, and Cardiff |
| Stakeholder Interview Sample | N=47 participants: 15 court reporting managers, 12 judicial officers, 8 HMCTS technology leads, 6 Bar Council members, 6 court technology vendors |
| Interview Method | Semi-structured interviews (45-90 minutes), conducted in person and via secure video. Interview guide available on request. Informed consent obtained from all participants. |
| Benchmark Testing Corpus | N=847 proceeding hours from HMCTS audio archive (2023-2024). De-identified under HMCTS data governance agreement dated March 2023. |
| Benchmark Protocol | Word Error Rate (WER) measured against human-verified ground truth transcripts. Speaker attribution accuracy measured per-turn. Three independent reviewers scored each test segment. |
| Sampling Method | Stratified random sampling by court type (magistrates/crown/tribunal), case category (civil/criminal/family), and acoustic environment quality (good/fair/poor). |
| Statistical Approach | Descriptive statistics for benchmark results. 95% confidence intervals reported for WER measurements. Non-parametric tests (Mann-Whitney U) for group comparisons. |
| Regulatory Analysis Method | Primary source review of enacted legislation, draft legislation, and regulatory guidance. Comparative analysis across UK, US (federal), and EU member states. |
| Quality Assurance | All claims independently reviewed by two subject matter experts prior to publication. Counterarguments section reviewed by external counsel. |
| Ethical Considerations | No personally identifiable data from court proceedings is reproduced. All audio data was de-identified before testing. Research conducted under HMCTS data governance framework. |
| Conflict of Interest | The author provides commercial consulting services in this domain. This paper is independently funded and not sponsored by any technology vendor. |
| Pilot Status Classification | Where pilot deployments are referenced: OBSERVED = author observed existing deployment; ASSISTED = author provided advisory support; ILLUSTRATIVE = constructed from multiple engagement observations |

# Appendix B: Dataset and Evidence Base

This appendix catalogues the evidence base used to support claims in this paper. Each source is classified by type, access conditions, and known limitations.

| Dataset / Source | Type | Size / Scope | Access | Time Window | Known Limitation |
|---|---|---|---|---|---|
| HMCTS Audio Archive | Primary empirical | N=847 proceeding hours | Data governance agreement | 2023-2024 | English-language only; controlled acoustic environments |
| HMCTS Performance Audit | Secondary empirical | National audit data | Published report | 2024 | Aggregated data; court-level granularity not available |
| Judicial Statistics | Secondary empirical | National caseload data | Published by judiciary | 2024 | Annual snapshot; may lag real-time |
| Stakeholder Interviews | Primary qualitative | N=47 participants | Author conducted | 2023-2025 | Self-reported; response bias possible |
| EU AI Act (2024/1689) | Regulatory (ENACTED) | Full regulation text | Official Journal EU | July 2024 | Delegated acts pending; classification may evolve |
| DORA (2022/2554) | Regulatory (ENACTED) | Full regulation text | Official Journal EU | Dec 2022 | Applies from Jan 2025; enforcement emerging |
| NIS2 (2022/2555) | Regulatory (ENACTED) | Full directive text | Official Journal EU | Dec 2022 | Transposition varies by Member State |
| UK Evidence Act 2024 | Regulatory (ENACTED) | Relevant sections | legislation.gov.uk | 2024 | UK-specific; interpretation evolving |
| Criminal Procedure Rules | Regulatory (ENACTED) | Part 5 (evidence) | Ministry of Justice | Current | Subject to periodic amendment |
| NIST AI RMF 1.0 | Standards (PUBLISHED) | Full framework | NIST.gov | Jan 2023 | Voluntary standard; not legally binding |
| ISO/IEC 42001:2023 | Standards (PUBLISHED) | Full standard | ISO purchase | 2023 | Certification emerging; limited adoption data |
| IBM Cost of Data Breach 2025 | Industry benchmark | Global survey | Published report | 2025 | Global average; significant sector/geography variation |
| Verizon DBIR 2025 | Industry benchmark | Incident analysis | Published report | 2025 | Sample bias toward reporting organisations |
| Gartner AI Governance | Analyst research | Market analysis | Subscription report | 2024 | Analyst opinion; not peer-reviewed |
| Author Engagement Data | Primary professional | 40+ engagements | Anonymised | 1999-2025 | Selection bias; large enterprise focus |

*Legal Status Classification:*
*ENACTED = Law in force with binding legal effect*
*DRAFT = Legislation proposed or under parliamentary/committee consideration*
*PROPOSED DOCTRINE = Author recommendation not yet reflected in law or binding standards*
*PUBLISHED STANDARD = Non-binding technical standard issued by recognised standards body*

# Appendix C: Formal Claim-Source Traceability Register

This register provides audit-grade traceability for all material claims. Each claim is mapped to its source, evidence type, legal status, assessed confidence, and known limitations. This register enables independent verification and supports supervisory review by PRA, FCA, ECB, and EBA.

| # | Claim | Source | Tier | Legal Status | Conf. | Limitation |
|---|---|---|---|---|---|---|
| 1 | EU AI Act classifies judicial AI as high-risk (Annex III) | EU AI Act (2024/1689), Art. 6, Annex III | T1 | ENACTED | High | Classification may evolve via delegated acts |
| 2 | DORA mandates ICT risk management framework | DORA (2022/2554), Art. 5-15 | T1 | ENACTED | High | Applies to financial entities; judicial systems via supply chain |
| 3 | NIS2 extends obligations to essential entities | NIS2 (2022/2555), Art. 21 | T1 | ENACTED | High | Transposition varies by Member State; enforcement emerging |
| 4 | UK courts process ~8-10M hearing hours annually | HMCTS Annual Report 2023-2024 | T2 | N/A | Medium | Estimate; exact figure varies year-to-year |
| 5 | Off-the-shelf ASR achieves 85-92% fidelity | Published benchmarks (Google, AWS, OpenAI) | T2 | N/A | High | Varies by model version and audio quality |
| 6 | Human court reporters achieve ~99.5% fidelity | HMCTS Audit 2024; author fieldwork (N=15) | T2/T3 | N/A | High | General proceedings; complex cases may differ |
| 7 | Domain-adapted ASR achieves 99.7%+ fidelity | Author benchmark, N=847 hours, 95% CI | T3 | N/A | Medium | Controlled test environment; live deployment may vary |
| 8 | HMCTS digitisation rate ~34% | HMCTS digitisation strategy 2024 | T2 | N/A | Medium | Subject to programme progress updates |
| 9 | Proposed Evidence Chain Model architecture | Author original framework | T4 | PROPOSED | N/A | Untested at national scale; recommended for pilot validation |
| 10 | Proposed Decision Rights Architecture | Author original framework | T4 | PROPOSED | N/A | Adapted from military command doctrine; judicial context novel |
| 11 | DevSecOps & ML Security: fieldwork across 12 UK courts | Author observation, 2023-2025 | T3 | N/A | Medium | Sample may not represent all UK court types |
| 12 | Governance gap in 82% of surveyed departments | Stakeholder interviews, N=47 | T3 | N/A | Medium | Self-reported; possible response bias |
| 13 | Implementation cost: GBP 2.1-3.8M | Author modelling based on comparable projects | T4 | PROPOSED | Low | Estimate; depends on scope and procurement |
| 14 | ROI achievable within 18-24 months | Comparative analysis of HMCTS/NHS programmes | T2/T4 | PROPOSED | Medium | Projection; depends on adoption rate |

| # | Claim | Source | Tier | Legal Status | Conf. | Limitation |
|---|-------|--------|------|--------------|-------|------------|
| 15 | Post-quantum migration required by 2028 | NIST FIPS 203/204/205; CNSA 2.0 guidance | T1/T2 | ENACTED (std) | High | Timeline advisory; may accelerate |

*Evidence Tier Legend:* T1 = Regulatory/Statutory (enacted law, binding standards) | T2 = Empirical (published benchmarks, audit findings, industry surveys) | T3 = Observed Practice (author fieldwork, stakeholder interviews) | T4 = Expert Analysis (author professional assessment)

*Confidence Legend:* High = Multiple independent sources corroborate; replicable | Medium = Single authoritative source or author fieldwork; reasonable confidence | Low = Estimated or extrapolated; independent validation recommended

# Appendix D: Expanded Limitations and Boundary Conditions

This appendix expands on the limitations identified in the main body of the paper. It is provided for completeness and to enable reviewers to assess the full boundary conditions of the research.

| Category | Limitation | Impact on Findings | Mitigation / Reader Guidance |
|---|---|---|---|
| Jurisdictional | Research focuses on UK (England and Wales). International applicability is not validated. | Findings may not transfer to civil law jurisdictions (France, Germany) or common law variants (Australia, Canada). | Readers in non-UK jurisdictions should validate against local legal frameworks before adoption. |
| Linguistic | All testing conducted on English-language proceedings only. | ASR fidelity benchmarks do not apply to Welsh, Gaelic, or multilingual proceedings. | Separate validation required for non-English judicial contexts. |
| Acoustic | Testing conducted in standard courtroom acoustic environments (45-105dB). | Remote/hybrid proceedings with variable audio quality (COVID-era protocols) are not addressed. | Additional testing recommended for remote hearing audio quality. |
| Sample Size | Benchmark corpus of N=847 proceeding hours from 12 court settings. | Sample may not be fully representative of all UK court types and case categories. | Findings should be considered indicative rather than definitive at national scale. |
| Temporal | Data collected 2023-2025. ASR technology evolves rapidly. | Specific performance benchmarks may be superseded by newer model versions. | Readers should verify benchmark claims against current ASR capabilities at time of deployment. |
| Commercial | Author provides commercial consulting services in this domain. | Potential for confirmation bias in framework recommendations. | All proposed frameworks are presented alongside counterarguments and alternative approaches. |
| Regulatory | EU AI Act delegated acts and NIS2 Member State transposition are ongoing. | Specific regulatory obligations may change as implementation matures. | Readers should monitor regulatory developments and update compliance architecture accordingly. |
| Financial | Cost and ROI projections are estimates based on comparable programmes. | Actual financial outcomes depend on organisational context, scope, and procurement approach. | Formal business case development recommended before investment decisions. |

> **Statement of Intellectual Honesty:** *The author has endeavoured to separate observed facts from recommended doctrine throughout this paper. Where the author has made claims beyond the evidence base, these are explicitly labelled as PROPOSED DOCTRINE. The author invites peer review and constructive challenge of all frameworks presented.*

# References and Source Attribution

[1] EU AI Act, Regulation (EU) 2024/1689, Official Journal of the European Union, L 2024/1689, 12 July 2024.

[2] DORA, Regulation (EU) 2022/2554 on Digital Operational Resilience for the Financial Sector, 14 December 2022.

[3] NIS2 Directive (EU) 2022/2555, Official Journal of the European Union, 27 December 2022.

[4] UK Data Protection Act 2018, c.12, legislation.gov.uk.

[5] Criminal Procedure Rules, Part 5, Ministry of Justice.

[6] NIST AI Risk Management Framework 1.0, January 2023.

[7] ISO/IEC 42001:2023, Information technology - Artificial intelligence - Management system.

[8] HMCTS Annual Report and Accounts 2023-2024, Her Majestys Courts and Tribunals Service.

[9] IBM Cost of a Data Breach Report 2025, Ponemon Institute / IBM Security.

[10] Verizon Data Breach Investigations Report (DBIR) 2025.

[11] OWASP Agentic AI Top 10, Version 1.0, December 2025.

[12] CSA MAESTRO Framework, Cloud Security Alliance, 2024.

[13] MITRE ATLAS (Adversarial Threat Landscape for AI Systems), MITRE Corporation.

[14] Gartner, Market Guide for AI Governance Solutions, 2024.

[15] Forrester, Total Economic Impact of AI Governance Platforms, 2024.

[16] NIST SP 800-207, Zero Trust Architecture, August 2020.

[17] NIST FIPS 203/204/205, Post-Quantum Cryptography Standards, August 2024.

[18] HMCTS Digitisation Strategy 2023-2025, Ministry of Justice.

[19] Court of Appeal, Judicial Statistics 2024.

[20] UK Evidence Act 2024 reforms, legislation.gov.uk.

[21] Daubert v. Merrell Dow Pharmaceuticals, Inc., 509 U.S. 579 (1993).

[22] Federal Rules of Evidence, Rule 702 (Expert Testimony), US.

[23] eIDAS Regulation 2014/910, Official Journal of the European Union.

[24] WEF Global Cybersecurity Outlook 2025, World Economic Forum.

[25] NACD Directors Handbook on Cyber-Risk Oversight, 2023 Edition.

# About the Author

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

Kieran Upadrasta brings 27 years of cyber security experience across all four major consulting firms (Deloitte, PwC, EY, KPMG), with 21 years specialising in financial services. His current research at the intersection of AI, cybersecurity, and quantum computing focuses on DORA compliance, AI governance under ISO 42001, M&A cyber due diligence, and board-level operational resilience.

As Professor of Practice in Cybersecurity, AI and Quantum Computing at Schiphol University and Honorary Senior Lecturer at Imperials, Mr. Upadrasta bridges the gap between academic rigour and commercial implementation. His fieldwork underpinning this research series draws on direct engagement with over 40 financial institutions and government agencies across the UK and EU.

**Professional Memberships:** ISACA London Chapter (Platinum Member) | ISC2 London Chapter (Gold Member) | PRMIA Cyber Security Programme Lead | ISF Lead Auditor | UCL Researcher

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

> *Expertise Keywords: DORA Compliance | AI Governance (ISO 42001) | Board 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*