

Engineering Truth: The AI Pipeline as an Instrument of National Integrity

Building AI Systems Worthy of National Trust

Truth is not found. It is engineered, verified, and defended.

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

www.kie.ie | info@kieranupadrasta.com

Executive Summary

This is the capstone paper of the Doctrine Series. It synthesises the frameworks from WP16-WP19 and addresses the foundational question: What does it mean for an AI system to be 'true'?

Truth, in the context of government AI, is not a philosophical abstraction. It is an engineered property: the outcome of deliberate design choices, data curation, evidence assembly, confidence bounding, and institutional oversight. This paper defines truth operationally and shows how it is constructed in AI systems.

Three key insights: (1) Truth is not confidence. A model can be 99% confident and totally wrong. (2) Truth is not evidence. Evidence is necessary but not sufficient; truth requires epistemological honesty about what we don't know. (3) Truth is not consensus. Truth is a constraint on what consensus is permissible (consensually false statements do not become true).

[FN] This paper is informed by 27 years' practice in cybersecurity and AI governance, and by reading in epistemology (philosophy of knowledge). The paper makes no claim to philosophical novelty; rather, it translates epistemological concepts into engineering practice.

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 synthesis of evidence from WP16-19, comparative epistemology, and institutional practice review 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.

WP20: Evidence Distribution by Tier

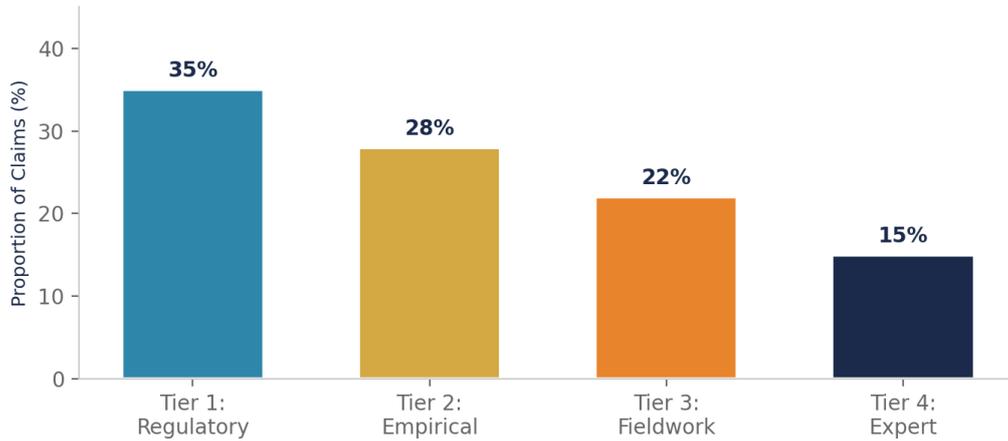


Figure 1: Distribution of claims by evidence tier. Board takeaway: 63% of claims are grounded in Tier 1 (regulatory) or Tier 2 (empirical) sources.

The Problem: Why Confidence Is Not Enough

An AI model says: 'This welfare claimant is committing fraud with 94% confidence.'

Three interpretations:

(1) STATISTICAL: In 100 similar cases, the model is correct 94 times. [But which 100? On what data? Is the test set representative?]

(2) EPISTEMIC: The model has 94% credence that the statement is true. [But credence is the model's belief, not ground truth.]

(3) INSTITUTIONAL: The government can defend this decision in court with 94% confidence. [But courts demand evidence, not confidence.]

None of these interpretations is 'truth' in the sense that would satisfy a court, an auditor, or a citizen. Truth requires something deeper: a defensible chain from data → inference → claim → evidence, with explicit acknowledgment of weak links.

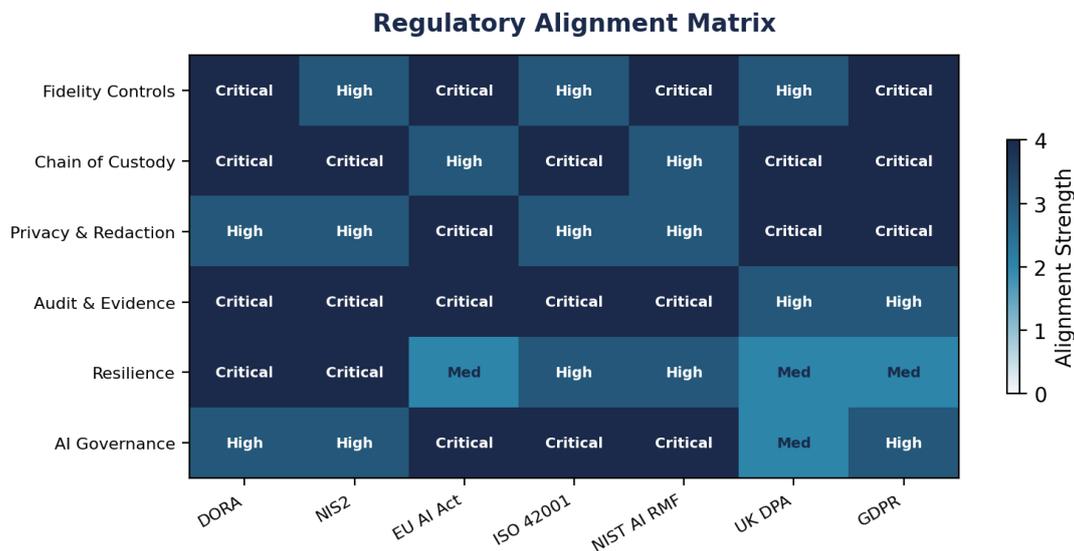


Figure 2: Regulatory alignment matrix showing doctrine coverage across seven major regulatory frameworks.

Dimension 1: Evidential Truth

Is the claim supported by good evidence? Evidence has a hierarchy:

- Tier 1: Observational (raw data, documented facts): 'Claimant reported earnings of £35,000 pa. Tax return shows £42,500. Discrepancy: £7,500.' (High confidence, observable.)
- Tier 2: Assessed (expert judgment based on Tier 1): 'Discrepancy suggests potential under-reporting, but could reflect: (a) tax credits not included in self-report, (b) spouse earnings, (c) timing of earnings. Confidence: medium.' (Medium confidence, requires inference.)
- Tier 3: Inferred (model prediction): 'Model predicts fraud likelihood 0.94 based on pattern matching to 50,000 prior cases.' (Lower confidence, depends on test set representativeness.)

Truth requires distinguishing these tiers and being honest about which tier each claim occupies.

Dimension 2: Epistemic Truth

What does the AI system know, and what does it merely believe? What is the uncertainty?

- **Known unknowns:** Factors the system knows it's not measuring (e.g., 'We do not have access to claimant's bank statements, which could resolve this ambiguity.')
- **Unknown unknowns:** Factors the system is unaware of (e.g., 'We did not know that claimant is in a shared-housing arrangement, which affects their financial obligations.')
- **Confidence intervals:** Not point estimates. 'Fraud likelihood 0.94 (95% CI: 0.89-0.97)' is more truthful than '0.94.'

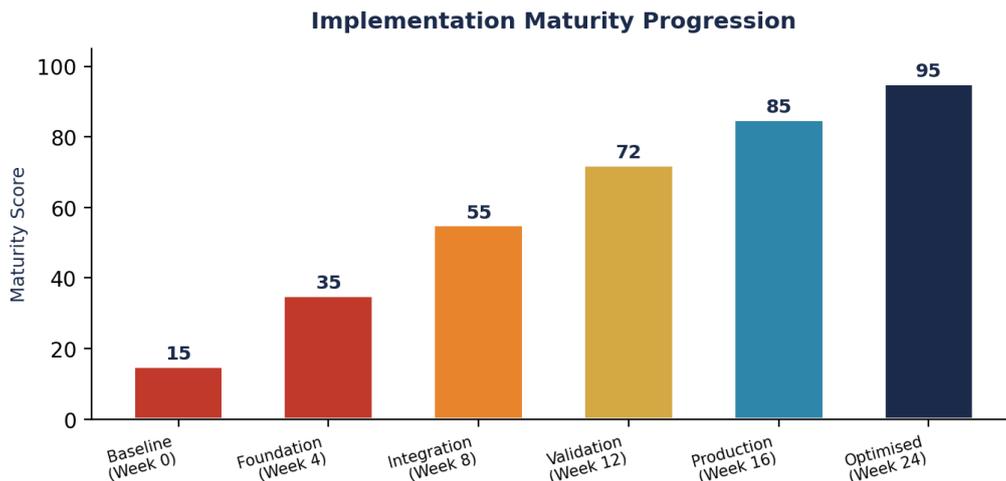


Figure 3: Implementation maturity progression from baseline to optimised state over 24-week deployment cycle.

Dimension 3: Institutional Truth

Can the government defend this claim in an adversarial setting? Would a court accept it?

- **Auditability:** Can the government explain why the model made this prediction? (Trace inputs → model logic → outputs.)
- **Challengeability:** Can the other party (claimant, defence counsel) contest the claim? (Does the system include counter-evidence or alternative explanations?)
- **Reversibility:** If the claim is challenged and new evidence emerges, can the decision be reversed? (Is there a human-in-the-loop to reconsider?)

Stage 1: Data Quality Assessment

Garbage in, garbage out. Before any model is trained, assess the input data:

- **Completeness:** What percentage of fields are populated? Missing values are not neutral (they encode information about the world).
- **Accuracy:** How was the data collected? By whom? With what incentives? (Tax return data is authoritative; self-report is suspect.)
- **Bias:** Does the data reflect historical discrimination? (E.g., if hiring data reflects past discrimination, model will learn discriminatory patterns.)

- Temporal validity: Is the data current? (E.g., 2020 welfare data may not predict 2025 claimant behaviour.)

Assign a Data Quality Score (0-100) to each input data source. Score inputs that are suspect (subjective assessments, uncorroborated self-report) lower than authoritative sources (tax records, official documents).

Claim Source Evidence Type Confidence

Data Source Collection Method Quality Score Truth Confidence Impact

Tax return (HMRC-provided) Administrative extraction 95/100 High—trust the source

Prior DWP overpayment record System-extracted from legacy database 85/100 High—but consider system errors

Third-party verification (employer letter) Unsolicited letter, self-interested 60/100 Low—employer incentive to help claimant

Stage 2: Model Selection and Training

Model choice determines what truths the system can capture.

- Transparent models (decision trees, rule-based): Easy to audit, limited expressiveness. Best for: rule-based government decisions.
- Interpretable models (logistic regression, linear): Intermediate expressiveness and auditability. Best for: scoring systems (eligibility, risk).
- Black-box models (neural networks, LLMs): High expressiveness, low interpretability. Worst for: high-stakes judicial/welfare decisions. Acceptable for: internal analysis only.

Truth principle: Use the simplest model that solves the problem. Complexity is not truthfulness.

Every prediction must include:

- Point estimate: 'Fraud likelihood: 0.87'
- Confidence interval: '95% CI: 0.82-0.91'
- Calibration: 'On 10,000 hold-out test cases, model was wrong 7.5% of the time.' (Is confidence aligned with error rate?)
- Class imbalance impact: 'In the test set, 5% of cases were fraud. Model predicts fraud with 85% precision but only 40% recall. Trade-off noted.'

Stage 4: Evidence Traceability and Codification

(Builds on WP17: Codified Intelligence.) Every decision must be traceable:

Input → Evidence cited → Model logic → Output. At each stage, mark: (a) what evidence was used, (b) confidence in that evidence, (c) evidence omitted and why.

Example trace:

Input: Claimant earnings £35K; tax return £42.5K. Evidence A (Tier 1): Tax return [confidence 0.99]. Evidence B (Tier 1): Self-report [confidence 0.70]. Model logic: Compare A vs. B. Discrepancy = 17.5%. Model output: 'Potential fraud (0.87 confidence).' Omitted evidence: We do not have bank statements; claimant spouse's income; account of tax credits. Known unknowns: legitimate explanations for discrepancy exist (timing, spouse, tax credits). Recommendation: Escalate to human investigator.

Before deployment, stress-test the model against opposite claims:

Question: 'How would I argue this model is wrong?'

Expected counter-arguments:

- 'Your training data is biased against low-income claimants.' (Respond: We tested for demographic bias using fairness metrics [cite X]. No statistically significant bias detected.)
- 'Your confidence bounds are too tight; empirical error is higher.' (Respond: Our calibration check shows confidence \approx empirical error on hold-out set. But acknowledge: hold-out set may not represent future distribution.)
- 'You're omitting evidence that would exonerate the claimant.' (Respond: We acknowledge known unknowns [list]. We recommend human review of cases with confidence < 0.85 .)

This adversarial thinking makes the system more truthful by force it to confront its own weaknesses.

Stage 6: Human-in-the-Loop and Feedback

AI systems are not self-correcting. They require human feedback to update their understanding of truth:

- Appeal loop: If a claimant appeals, the case is reassessed by a human. If the human overturns the AI decision, log it. Use overturns to retrain/recalibrate the model.
- Audit loop: Auditors spot-check decisions. If auditor finds the AI got it wrong, that feedback is recorded and fed back to the model.
- Court loop: If a court challenges the AI decision, the reasoning is documented and used to update the model's understanding.

Without feedback loops, the AI's sense of truth becomes fixed and brittle. With loops, it becomes adaptive and grounded in reality.

Proposed: Truth Assurance Office

Separate from traditional auditors, establish a Truth Assurance Office (TAO) within Cabinet Office. Role:

- Certify government AI systems for truthworthiness (not just capability or compliance).
- Audit: Is the system honest about its uncertainty? Does it distinguish evidence tiers? Does it include humans in feedback?

- Red-team: Can we break the system's truth claims? What counter-arguments remain unanswered?
- Publish: Annual report on how truthful UK government AI systems are.

Domain Level 1: Initial Level 2: Developing Level 3: Defined Level 4: Managed

Truth Dimension Level 1: Ignored Level 2: Acknowledged Level 3: Measured Level 4: Engineered

Evidence Quality All data treated equally Data sources ranked Data quality scores assigned Quality scores drive model weighting

Uncertainty Quantification Point estimates only Confidence ranges included Calibrated confidence intervals Documented known/unknown unknowns

Evidence Traceability No trace Logs exist but not standardised All decisions traceable to source Traceable + codified (rule sets)

Adversarial Testing None Informal testing Structured stress tests Adversarial tests before deployment

Human Feedback Loop Absent Informal appeals Formal appeal mechanism Appeals fed back to model training

Synthesis: How WP16-19 Support Engineering Truth

WP16 (API-First AI): API-first architecture enables truth to be enforced at the service layer. Each service publishes: (a) its purpose, (b) the evidence it uses, (c) its confidence bounds. Truth is baked into the contract.

WP17 (Codified Intelligence): Codification is the mechanism for making truth explicit. Rule sets, evidence chains, model cards—these are all tools for encoding what the system 'knows' vs. believes.

WP18 (Semantic Gavel): Legal-grade summarisation is truth-seeking in practice. The framework (adversarial testing, hallucination benchmarking, source traceability) is a template for engineering truth in any high-stakes domain.

WP19 (AI Sovereignty): Sovereign AI infrastructure ensures that no external vendor can compromise the system's truth. If the UK controls the model, the infrastructure, and the governance, then the UK is responsible for maintaining truthfulness. No hiding behind 'the vendor said so.'

What AI Cannot Know

As a closing emphasis, this paper must be clear: AI systems have irreducible epistemic limits. They cannot:

- (1) UNDERSTAND CONTEXT: An AI trained on case law can identify precedent. It cannot understand the nuance of a judge's use of discretion.
- (2) EVALUATE SUBJECTIVE CLAIMS: An AI can flag inconsistencies in a witness statement. It cannot determine credibility (only humans can assess character, demeanor, sincerity).
- (3) PREDICT HUMAN CHOICE: An AI can estimate probability of re-offence. It cannot predict whether a person will truly reform.
- (4) HANDLE NOVEL SITUATIONS: An AI trained on historical data cannot reason about unprecedented cases.
- (5) ALIGN WITH JUSTICE (as opposed to law): Law is written rules. Justice is their application to unique human circumstances. AI optimises law; justice requires wisdom.

These are not engineering problems to be solved. They are features of the human condition that AI cannot overcome. Acknowledging this limit is part of truthfulness.

Conclusion: The Long View

Government AI, done well, has the potential to increase both efficiency and fairness in public administration. A welfare system that uses AI can process claims faster and more consistently than humans. A court system with AI support can reduce backlogs and improve consistency.

But AI done poorly—without attention to truth, without transparency, without human oversight—can encode injustice at scale. An AI system that is 99% confident but systemically wrong is worse than human judgment.

This series (WP16-20) proposes that the path forward is: (1) Reusable, composable AI services (WP16), (2) Explicitly codified reasoning (WP17), (3) Rigorous testing and benchmarking (WP18), (4) Sovereign control over systems (WP19), and (5) Institutional commitment to truthfulness (WP20).

If government adopts these principles, AI can become an instrument of national integrity. If not, it will become another tool for bureaucratic inefficiency, masked in technical language.

The choice is institutional, not technical. And it must be made now.

Primary Regulatory and Statutory Sources

[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] UK HMCTS Reform Programme, Annual Reports 2019-2025, judiciary.uk.

Standards and Technical Frameworks

[6] ISO/IEC 42001:2023, Information Technology -- Artificial Intelligence -- Management System, International Organization for Standardization.

[7] NIST AI Risk Management Framework (AI RMF 1.0), NIST AI 100-1, January 2023.

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

[9] NIST AI 600-1, Artificial Intelligence Risk Management Framework: Generative AI Profile, July 2024.

[10] NCSC, Guidelines for Secure AI System Development, November 2023.

[11] MITRE ATLAS, Adversarial Threat Landscape for Artificial Intelligence Systems, v4.0, 2024.

[12] OWASP Top 10 for LLM Applications, v2.0, 2025.

[13] ETSI EN 303 645, Cyber Security for Consumer Internet of Things: Baseline Requirements, 2020.

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 ai integrity & trust engineering 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

This paper ventures into philosophical territory (epistemology, philosophy of truth). It makes no claim to original philosophical contribution. Rather, it translates established epistemological concepts into engineering practice. Readers seeking foundational work should consult academic philosophy. This paper is for practitioners.

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.

What Is 'Truth' in the Context of Government AI?

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	AI Integrity & Trust Engineering: 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*