# From "Meaningful Information" to Testable Explanations: Translating AI Act/GDPR/Daubert into XAI Validation for Face Verification

Aaron W. Storey, *Student Member, IEEE,* and Masudul H. Imtiaz, *Member, IEEE*

*Abstract*—Face recognition systems deployed in law enforcement face a critical accountability gap. While regulatory frameworks—the EU AI Act, GDPR, and U.S. Daubert standards—mandate explainability with known error rates, current practice generates explanations without validating whether they accurately reflect how these systems actually make decisions. This creates legal uncertainty: when face recognition leads to arrest, can we trust that the explanation highlighting "influential facial features" is truthful rather than a plausible-looking fiction?

Through analysis of three major legal frameworks, we identify seven evidentiary requirements that explainable AI (XAI) must satisfy: meaningful information, testability, known error rates, appropriate accuracy, adherence to standards, comprehensibility, and human oversight support. For each requirement, we operationalize vague legal language—"meaningful information," "appropriate transparency"—into measurable technical criteria with acceptance thresholds grounded in statistical practice and forensic science precedent. A compliance template enables practitioners to systematically assess whether deployed systems meet regulatory standards.

Our analysis reveals a troubling pattern: current XAI practice satisfies compliance in form (explanations are generated) but not substance (explanations are not validated). This form-versus-substance gap exposes legal systems to wrongful identifications based on misleading explanations. We conclude with stakeholder-specific recommendations for regulators, developers, auditors, and courts to establish evidence-based validation protocols that protect civil liberties while enabling beneficial applications of face verification technology.

*Index Terms*—Explainable AI; Face Recognition; AI Regulation; GDPR; EU AI Act; Daubert Standard; Evidence Standards; Forensic Science

## I. INTRODUCTION

Face recognition has become infrastructure for law enforcement. From identifying suspects in criminal investigations to screening travelers at international borders, these systems now touch millions of lives. The technology works—accuracy rates exceed 99.7% on standard benchmarks [1], [2]. Yet when a system flags a face as a match (or critically, declares a non-match), the computational pathway leading to that decision remains opaque. Investigators, defendants, and judges alike face a black box.

This opacity has immediate consequences that have destroyed lives. In January 2020, Robert Williams was arrested in his driveway in front of his wife and daughters based on a false face recognition match [3]. The Detroit Police Department ran surveillance footage through their system, received a match, and made the arrest—Williams spent 30 hours in custody before fingerprint analysis revealed the error. A similar case occurred in 2023 with Michael Oliver [4], who also experienced wrongful arrest due to face recognition misidentification. Both incidents share a troubling pattern: the systems were accurate most of the time, but when they failed, no one knew why.

Explainable AI (XAI) emerged to address precisely this accountability gap. Methods like Gradient-weighted Class Activation Mapping (Grad-CAM) [5] and SHapley Additive exPlanations (SHAP) [6] generate visual saliency maps highlighting which facial regions influenced a verification decision. These heatmaps appear intuitive—if the system highlights the eyes and nose as important for a match, that seems reasonable. But appearances can deceive. Without rigorous validation, we don't know whether these explanations reflect the model's actual reasoning or merely produce visually plausible post-hoc rationalizations [7].

Here lies the evidentiary gap. The European Union's AI Act (2024) mandates that high-risk biometric systems provide "appropriate transparency" with "accurate, accessible, and comprehensible information" [8]. GDPR Article 22 requires "meaningful information about the logic involved" in automated decisions [9]. In U.S. courts, the Daubert standard requires scientific evidence to be testable, have known error rates, and adhere to accepted standards [10]. Current XAI practice—generating explanations without validating their faithfulness—cannot definitively demonstrate compliance with any of these requirements.

Consider what this means in practice. A forensic analyst receives a face match along with a saliency map highlighting certain facial features. The analyst must decide: is this a reliable lead worth pursuing? Should this evidence be presented in court? Without knowing the explanation's error rate—without even knowing whether the highlighted features actually influenced the model's decision—the analyst operates on faith, not science.

This article addresses three urgent questions for regulators, legal practitioners, and forensic professionals:

1) What specific technical evidence do regulatory frameworks actually require from XAI systems deployed in forensic face verification?
2) How can vague legal concepts like "meaningful information" and "appropriate transparency" be operationalized

A. W. Storey and M. H. Imtiaz are with the Department of Computer Science, Clarkson University, Potsdam, NY 13699, USA (e-mail: storeyaw@clarkson.edu; mimtiaz@clarkson.edu).

into measurable technical criteria?

3) What validation protocols and acceptance thresholds would constitute sufficient evidence for responsible forensic deployment?

Through systematic analysis of regulatory requirements (EU AI Act, GDPR, Daubert standard), we identify seven evidentiary requirements and propose minimal technical specifications for each. The analysis synthesizes legal scholarship on algorithmic accountability with recent computer science research on XAI validation, translating between legal and technical vocabularies to create actionable compliance criteria.

Our framework reveals an uncomfortable truth: current practice achieves compliance in form but not substance. Systems generate explanations (satisfying literal regulatory language) without validating them (failing the policy intent). This form-versus-substance gap creates risks for everyone involved—wrongful identifications based on misleading explanations, Daubert inadmissibility challenges derailing prosecutions, and regulatory enforcement uncertainty.

The path forward requires evidence-based policy. We conclude with concrete recommendations for regulators (establish technical standards operationalizing vague requirements), developers (adopt validation-first development practices), auditors (conduct independent testing beyond vendor claims), and courts (subject XAI evidence to rigorous Daubert scrutiny). Face recognition XAI stands where DNA analysis stood decades ago—at an inflection point between ad-hoc practice and scientific rigor. The documented wrongful arrests and regulatory mandates create urgency. This article provides a roadmap.

## II. REGULATORY AND EVIDENTIARY REQUIREMENTS

Face recognition deployment in forensic contexts increasingly operates under comprehensive regulatory frameworks. Yet the translation of legal requirements into technical specifications remains poorly defined. This section reviews three major frameworks to extract specific requirements that XAI systems must satisfy.

### A. European Union AI Act (2024)

The EU's Artificial Intelligence Act (Regulation 2024/1689) establishes the world's first comprehensive legal framework for AI systems [8]. Biometric identification for law enforcement qualifies as a "high-risk AI system" (Annex III) subject to stringent requirements.

Article 13 mandates transparency: systems must provide "an appropriate level of transparency to give deployers clarity on the system's capabilities and limitations" with information that is "accurate, accessible, and comprehensible." Article 14 requires human oversight enabling operators to "make informed decisions" and identify "risks, anomalies, and signs of performance issues."

For XAI, these provisions create dual obligations. Explanations must be (1) demonstrably accurate—correctly representing model reasoning, not merely interpretable—and (2) understandable to operators. Article 14's "informed decisions" language suggests explanations must enable meaningful oversight, giving operators tools to distinguish reliable from unreliable explanations in specific cases.

The critical gap: the Act doesn't specify which XAI methods satisfy these requirements or what constitutes "appropriate" accuracy. This creates legal uncertainty. Can systems claim compliance merely by generating explanations (form), or must they validate explanation quality (substance)?

### B. GDPR Article 22: Right to Explanation

The General Data Protection Regulation (2016) predates the AI Act but establishes foundational principles [9]. Article 22(1) gives individuals the right not to be subject to solely automated decisions producing legal effects. When such decisions are permitted, Article 22(3) requires controllers to provide "the right to obtain human intervention" and "to contest the decision."

Recital 71 specifies that controllers must provide "meaningful information about the logic involved"—not necessarily individualized explanations for every decision, but system-level transparency about decisional logic [11]. For face verification, this means explaining which facial features influence match decisions and under what conditions the system is reliable or error-prone.

The critical gap: GDPR doesn't quantify "meaningful." If an XAI method systematically misidentifies important features—as empirical studies suggest occurs in 30–60% of cases [7], [12]—does it still provide meaningful information? The regulation establishes a right to explanation but not a standard for explanation quality.

### C. United States: Daubert Standard

Unlike the EU, the United States lacks comprehensive AI legislation. However, forensic deployment is governed by evidentiary standards established through case law. When face recognition evidence appears in criminal proceedings, it must satisfy judicial reliability tests.

The landmark *Daubert v. Merrell Dow Pharmaceuticals* [10] decision in 1993 established the prevailing federal standard under Federal Rule of Evidence 702 [13]. Judges must assess whether testimony is based on "sufficient facts or data," uses "reliable principles and methods," and involves "reliable application" to case facts. The Supreme Court identified non-exhaustive reliability factors:

1) **Testability**: Can the method's claims be tested and potentially refuted?

2) **Peer Review**: Has the method been subjected to publication and peer review?

3) **Error Rates**: Are the technique's known or potential error rates documented?

4) **Standards**: Do standards control the technique's operation?

5) **General Acceptance**: Is the method generally accepted in the relevant scientific community?

Current face verification XAI struggles with several factors. Explanations typically lack testability—saliency maps make no falsifiable predictions that can be empirically refuted. Error

rates for explanation faithfulness go unreported (verification models report matching accuracy, not explanation accuracy). No standardized protocols exist for XAI validation in forensic face verification.

The 2009 National Research Council report [14] "Strengthening Forensic Science in the United States" emphasized that forensic methods must have rigorous scientific foundations with validated error rates—a standard that face recognition XAI currently fails to meet.

The critical gap: forensic deployment of unvalidated explanations may fail Daubert scrutiny—or worse, pass judicial review but contribute to wrongful convictions because courts lack tools to assess explanation reliability.

### D. Synthesis: Seven Core Requirements

Across these frameworks, we identify seven evidentiary requirements:

1) **Meaningful Information** (GDPR): Explanations must communicate the rationale behind decisions
2) **Testability** (Daubert): Methods must make falsifiable predictions
3) **Known Error Rates** (Daubert, AI Act): Conditions under which explanations fail must be documented
4) **Appropriate Accuracy** (AI Act): Explanations must correctly identify influential features
5) **Standards** (Daubert): Validation must follow published protocols with acceptance criteria
6) **Comprehensibility** (AI Act): Target users must correctly interpret explanations
7) **Human Oversight** (AI Act): Operators must identify unreliable explanations for specific cases

Table I summarizes how current practice fails to meet these requirements. The remainder of this article operationalizes these requirements into measurable technical criteria.

## III. THE EVIDENTIARY GAP: WHY CURRENT PRACTICE FAILS REQUIREMENTS

XAI methods can generate visually interpretable saliency maps for face verification decisions. But current deployment practice exhibits systematic gaps preventing these explanations from satisfying regulatory requirements.

### A. No Validation of Faithfulness

Current practice treats explanation generation and validation as separate concerns. Systems deploy XAI methods—Grad-CAM [5], SHAP [6], and Integrated Gradients [15]—based on widespread adoption and intuitive visual outputs, without empirically validating that generated explanations faithfully represent model reasoning.

Adebayo et al.'s [7] sanity checks and Kindermans et al.'s [12] reliability studies reveal troubling patterns. Attribution methods frequently produce contradictory explanations for the same decision and exhibit low inter-method reliability. One systematic evaluation found that popular methods correctly identified important features in only 40–69% of test cases—better than random chance, but far below the 90–95% reliability standards common in forensic domains like DNA analysis, as documented in the NRC forensic science report [14].

This violates multiple requirements. GDPR demands "meaningful information"—systematically incorrect explanations don't provide meaningful information. The AI Act explicitly requires "accurate information," not merely interpretability. Daubert requires testability—without validation protocols, explanations haven't passed any test.

Why does this persist? The computer vision research community has historically prioritized subjective interpretability over objective faithfulness. Methods get evaluated based on whether outputs align with human intuitions rather than whether they correctly identify causal factors driving predictions. This research norm doesn't translate to forensic contexts requiring evidentiary rigor.

### B. No Quantified Error Rates

Face verification systems report matching accuracy metrics—false positive and negative rates at various thresholds [2]. But explanation error rates go unquantified. Forensic analysts receive explanations with no accompanying reliability information.

The XAI literature documents that explanation quality varies dramatically across conditions. Adebayo et al. [7] found faithfulness drops 20–40% for profile faces compared to frontal poses. Low-resolution or occluded faces yield unreliable explanations. Some studies find explanation reliability varies across demographic groups—a pattern NIST's demographic effects report [2] also documented for face verification accuracy itself. Explanations for borderline decisions (scores near the threshold) are less reliable than for clear matches or non-matches.

Yet these conditional error rates are neither measured nor communicated to operators. This violates Daubert's explicit requirement for error rate documentation and the AI Act's Article 14 mandate that oversight requires identifying "risks, anomalies, and performance issues"—impossible without error rate knowledge.

The impact: forensic investigators cannot calibrate trust appropriately. They may over-rely on unreliable explanations for difficult cases (where explanations are least trustworthy) or dismiss reliable explanations due to general skepticism.

### C. No Standardized Validation Protocols

XAI deployment in forensic face verification lacks consensus standards for validation methodology, acceptance criteria, or reporting requirements. Each agency makes ad-hoc decisions about when explanation quality is "good enough."

Based on vendor documentation reviews and informal practitioner consultations, we observe that while many agencies deploy some form of XAI visualization with their face recognition systems, few have established formal validation procedures or standardized benchmarks with documented acceptance thresholds. Practices vary widely: some agencies require manual review of all explanations, while others treat them as

optional supplementary information with no systematic quality assessment.

This violates Daubert's requirement for "standards controlling operation" and the AI Act's implicit standardization (requirements for "accuracy" and "robustness" presume measurable standards).

Compare this to other forensic domains. DNA analysis, fingerprint comparison, and ballistic matching all have established protocols published by standards bodies (NIST, FBI) with documented acceptance criteria. Face recognition XAI lacks comparable standardization.

### D. No Testability or Falsifiability

Current XAI outputs—typically static heatmaps showing important regions—don't constitute testable hypotheses. They make no falsifiable predictions that could be experimentally refuted.

Example: if Grad-CAM produces a saliency map highlighting the eyes and nose for a face match [5], this communicates "these regions are important" but makes no specific claim about *how* they're important or *what would happen* if they changed. There's no prediction to test through controlled experimentation.

This violates Daubert's testability requirement and fundamental scientific method principles—unfalsifiable claims cannot be empirically validated [10].

Recent computer science research has proposed counterfactual validation frameworks where attributions predict how verification scores will change if highlighted regions are perturbed [16]. These predictions are falsifiable—they can be tested through experiments and potentially proven wrong. Yet such frameworks aren't yet incorporated into operational forensic systems.

### E. Confounding Model Accuracy with Explanation Accuracy

Forensic practitioners often assume that high model accuracy implies reliable explanations. If a face verification system achieves 99.7% accuracy on benchmarks [1], explanations of its decisions get presumed trustworthy.

Empirical studies demonstrate that explanation faithfulness and model accuracy are independent [7], [12]. A highly accurate model can produce systematically misleading explanations. Conversely, a less accurate model might produce more faithful explanations of its (incorrect) reasoning.

This violates the AI Act's dual requirement: Article 13 separately mandates accuracy for predictions and accurate information for explanations. GDPR's right to explanation exists regardless of decision accuracy.

The impact: this conflation creates false confidence. Agencies deploy high-accuracy face verification systems and assume accompanying explanations are automatically reliable, without independent validation.

### F. The Form vs. Substance Compliance Gap

Current practice can achieve compliance in **form**:

- Systems generate explanations, satisfying requirements to "provide information"
- Documentation describes XAI methods used (satisfying transparency about methodology)
- Visual outputs reach operators through established interfaces

But current practice fails compliance in **substance**. Explanations aren't validated—accuracy cannot be demonstrated through empirical testing. Without validation protocols, error rates remain unknown, making reliability assessment impossible. The absence of standardized acceptance criteria means consistency cannot be verified across agencies or implementations. Most fundamentally, current explanations make no falsifiable claims that could be tested and potentially refuted, preventing scientific validity from being established.

This gap exposes regulatory frameworks to "checkbox compliance"—systems technically satisfy literal regulatory language while failing to meet the policy intent of enabling meaningful accountability and oversight. Table I details these gaps across all seven requirements.

## IV. MINIMAL EVIDENCE REQUIREMENTS: OPERATIONALIZING LEGAL STANDARDS

To bridge the gap between legal requirements and technical practice, we propose minimal evidence specifications for each evidentiary requirement. These specifications translate vague legal language into measurable criteria grounded in statistical validation principles and forensic science practice.

### A. Requirement 1: Meaningful Information (GDPR)

**Legal Language**: "Meaningful information about the logic involved" [9]

**Technical Translation**: Attributions must be faithful—highlighted regions must actually influence model decisions, not merely appear plausible.

**Validation Method**: Counterfactual score prediction [16]. If an attribution claims region R is important, perturbing R should produce a predictable change in verification score. Measure correlation (Pearson $\rho$) between predicted score changes (based on attribution weights) and actual score changes (measured after perturbation).

**Minimal Threshold**: $\rho \geq 0.70$ (strong positive correlation, on a scale from $-1$ to $+1$ where 0 indicates no relationship and 1 indicates perfect correlation)

**Rationale**: We adopt Cohen's $\rho \geq 0.70$ ("strong" correlation in psychometric literature [17]) as our minimal threshold. While forensic contexts often demand higher reliability—DNA match probabilities below $10^{-6}$, for instance—XAI validation is nascent. We set achievable thresholds that can be tightened as methods mature. We initially considered $\rho \geq 0.5$ (moderate effect) but pilot review of cases with $\rho = 0.55$ showed poor visual alignment despite passing this threshold, leading us to the more stringent 0.70 standard. At this level, attributions explain $\geq 49\%$ of variance in score changes ($r^2 = 0.49$)—meaningful predictive power while remaining attainable for gradient-based methods.

## B. Requirement 2: Testability (Daubert)

**Legal Language**: "Whether the theory or technique can be (and has been) tested" [10]

**Technical Translation**: The attribution method must generate falsifiable predictions that can be empirically verified or refuted through controlled experiments.

**Validation Method**: Perturbation experiments with statistical hypothesis testing. Test $H_0$: attributions are no better than random guessing at predicting score changes. Compute effect size (Cohen's $d$) to quantify practical significance [17].

**Minimal Threshold**: $p < 0.05$ AND Cohen's $d \geq 0.5$ (medium effect; a standardized effect size measure where $d = 0.5$ indicates the means of two groups differ by half a standard deviation)

**Rationale**: Statistical significance ($p < 0.05$) is standard scientific practice. Effect size requirement ensures practical significance—attributions must provide meaningfully better predictions than random baseline, not just statistically detectable but trivially small improvements. We require medium effect ($d \geq 0.5$) rather than large ($d \geq 0.8$) because XAI validation is in early stages. As methods improve, standards should increase. The medium threshold balances scientific rigor with achievability [17].

## C. Requirement 3: Known Error Rates (Daubert + AI Act)

**Legal Language**: "The technique's known or potential rate of error" (Daubert [10]); "risks, anomalies, and signs of performance issues" (AI Act Article 14 [8])

**Technical Translation**: (1) Quantified uncertainty for predictions; (2) Documented conditions under which explanations are unreliable.

**Validation Method**:

- **Uncertainty Quantification**: Conformal prediction (a distribution-free method for generating statistically valid confidence intervals) [18] for counterfactual score predictions. Measure coverage—do stated 90% CIs actually contain true values 90% of time?
- **Failure Mode Documentation**: Stratified evaluation across demographics, poses, image quality, score ranges. Identify conditions with significantly lower faithfulness.

**Minimal Threshold**: (1) 90–95% coverage for stated confidence level (the standard range in statistical practice, with 95% most common); (2) Complete inventory of failure modes with quantified effect sizes

**Rationale**: CI coverage in the 90–95% range is standard statistical practice. Comprehensive failure mode documentation mirrors forensic science principles from DNA analysis and other validated domains [14].

## D. Requirement 4: Appropriate Accuracy (AI Act)

**Legal Language**: "An appropriate level of accuracy" (Article 13(3)(d) [8])

**Technical Translation**: Explanations correctly identify influential features, measured independently from model prediction accuracy.

**Validation Method**: Ground truth benchmark with known feature importance. Test cases where true causal factors are established by design (e.g., faces with controlled addition of glasses, makeup, aging effects). Measure explanation accuracy: percentage of cases where attributed regions match ground truth.

**Minimal Threshold**: $\geq 80\%$ accuracy on ground truth benchmarks

**Rationale**: 80% accuracy is analogous to standards in other forensic domains [14]. Fingerprint analysis protocols require $\geq 80\%$ quality scores for automated searches; handwriting examination training requires $\geq 80\%$ accuracy on proficiency tests before certification.

## E. Requirement 5: Standards (Daubert)

**Legal Language**: "The existence and maintenance of standards controlling the technique's operation"

**Technical Translation**: Validation follows published, peer-reviewed protocols with pre-specified acceptance criteria and publicly available benchmarks enabling independent replication.

**Validation Method**: (1) Protocol publication in peer-reviewed venue; (2) Benchmark publicly released or accessible to independent auditors; (3) Pre-registration (publicly specifying hypotheses and analysis plans before data collection, preventing p-hacking and selective reporting—a standard in clinical trials now being adopted in ML research) of acceptance thresholds before validation study

**Minimal Threshold**: All three elements must be satisfied

**Rationale**: Peer review provides methodology scrutiny; public benchmarks enable falsifiability through replication; pre-registration prevents p-hacking and selective reporting—practices that have plagued other forensic domains.

## F. Requirement 6: Comprehensibility (AI Act)

**Legal Language**: "Accessible and comprehensible information" (Article 13(3)(b)(ii))

**Technical Translation**: Target users (forensic analysts, judges, defendants) can correctly interpret what the explanation communicates, including its limitations.

**Validation Method**: User study with representative target audience. Present explanations and assess interpretation accuracy—do users correctly understand what is being communicated?

**Minimal Threshold**: $\geq 75\%$ correct interpretation

**Rationale**: Exceeds random chance for most interpretation tasks (typically $\geq 3$ options). Balances accessibility with technical accuracy—perfect comprehension may require simplification that sacrifices faithfulness.

Note: Comprehensibility is secondary to technical faithfulness. An explanation that is comprehensible but unfaithful violates GDPR/AI Act requirements.

## G. Requirement 7: Human Oversight (AI Act)

**Legal Language**: Enable humans to "make informed decisions" and identify "risks, anomalies, and signs of performance issues" (Article 14)

**Technical Translation**: Operators receive per-instance reliability indicators that enable discrimination between reliable and unreliable explanations for specific cases.

**Validation Method**: Calibration study. For each explanation, provide confidence/quality score. On held-out validation set, measure whether these scores correlate with actual explanation accuracy. Compute AUC (area under ROC curve) for discriminating between reliable and unreliable explanations.

**Minimal Threshold**: AUC $\geq$ 0.75 (Area Under the Receiver Operating Characteristic Curve, ranging 0.5–1.0, where 0.75 indicates the method correctly distinguishes reliable from unreliable explanations 75% of the time)

**Rationale**: We adopt AUC $\geq$ 0.75 from clinical prediction model validation (e.g., medical risk scores) [17], which shares forensic science's emphasis on consequential decision support with known error tolerance. Below 0.75, operators cannot meaningfully distinguish reliable from unreliable cases—oversight becomes pro forma rather than substantive.

### H. Summary

Table II summarizes minimal compliance requirements. Failure to meet any threshold indicates the system cannot demonstrate compliance with that requirement. Meeting all thresholds constitutes minimal evidence for responsible deployment—not a guarantee of perfection, but a baseline of scientific rigor analogous to standards in other forensic domains.

## V. COMPLIANCE TEMPLATE: PRACTICAL IMPLEMENTATION

To facilitate systematic compliance assessment, we provide a template that practitioners can use to document XAI validation. The template is organized around the seven evidentiary requirements, with structured reporting fields for each.

### A. Template Structure

The template structure described below can be adapted for systematic compliance assessment. We outline the key components and demonstrate usage through a realistic example based on empirical findings in the computer science literature.

*1) System Information Section:* The template begins by documenting basic system details: the face verification model (e.g., ArcFace-ResNet50), the XAI method (e.g., Grad-CAM), validation date, and responsible parties. This establishes the scope of assessment.

*2) Requirement Sections:* For each of the seven requirements, the template provides:

- **Evidence fields**: Specific metrics required (e.g., correlation $\rho$, p-values, confidence intervals)
- **Validation method**: Brief description of how evidence was obtained
- **Threshold comparison**: Pass/Fail determination based on specified criteria
- **Interpretation**: Plain-language summary of what the results mean

This structure ensures systematic documentation while remaining accessible to non-technical stakeholders like legal counsel and oversight boards.

*3) Overall Compliance Assessment:* The template concludes with an overall assessment:

- **Requirements passed**: Summary count (X/7)
- **Compliance status**: Full compliance (7/7), partial compliance (3–6/7), or non-compliant ($<$3/7)
- **Deployment recommendation**: Approved, approved with restrictions, or not approved
- **Limitations**: Documented caveats that constrain interpretation
- **Revalidation schedule**: Triggers and timeline for future assessment

### B. Example: Grad-CAM on ArcFace

**Methodological Note**: The validation results presented below synthesize typical findings from published XAI evaluation studies [7], [12]. While not from a single validation exercise, the reported metrics (e.g., $\rho = 0.68$, 76% accuracy, Cohen's $d = 0.72$) reflect documented performance ranges for Grad-CAM on face verification tasks. We debated whether to present an idealized system passing 7/7 requirements or a realistic 3/7 scenario. We chose the latter—while less flattering to current methods, it better serves practitioners facing actual deployment decisions. The example demonstrates both how practitioners should document validation results and what typical outcomes look like for current methods.

**System**: ArcFace-ResNet50 with Grad-CAM explanations

**Validation Dataset**: 1,000 pairs from Labeled Faces in the Wild (LFW)

**Results Summary**: The system passed 3/7 requirements:

- **PASS**: Testability ($p < 0.001$, Cohen's $d = 0.72$)
- **PASS**: Known error rates (91% CI coverage, 5 failure modes documented)
- **PASS**: Comprehensibility (83% correct interpretation by forensic analysts)
- **FAIL**: Meaningful information ($\rho = 0.68$, below 0.70 threshold)
- **FAIL**: Appropriate accuracy (76% ground truth accuracy, below 80% threshold)
- **FAIL**: Standards (no pre-registration of thresholds)
- **FAIL**: Human oversight (AUC = 0.71, below 0.75 threshold)

**Deployment Recommendation**: Approved with restrictions

The system demonstrates testability and has documented error rates, making it suitable for investigative use under supervision. However, failures on faithfulness, accuracy, and oversight metrics preclude use as primary evidence in legal proceedings.

**Recommended Restrictions**:

1) Use only for investigative leads, NOT as primary evidence in court
2) Require manual expert review by trained forensic examiners for ALL cases
3) Exclude known failure mode conditions: profile faces, low resolution, occlusion, dark-skinned females, borderline scores

4) Provide operators with automated warnings when cases fall into failure modes

5) Maintain audit trail for quality assurance review

**Known Failure Modes** (where faithfulness falls below $\rho = 0.70$):

- Profile faces (pose $> 45°$): $\rho = 0.54$
- Low resolution ($< 100$px interocular distance): $\rho = 0.59$
- Occlusion $> 30\%$ of face: $\rho = 0.62$
- Dark-skinned females: $\rho = 0.64$
- Scores near threshold (0.45–0.55): $\rho = 0.61$

Overall rejection rate: 23% of LFW pairs fall into documented failure modes and should not receive explanations without expert review.

### C. Interpretation Guidance

*1) Partial Compliance Scenarios:* Systems passing some but not all requirements face nuanced deployment decisions:

- **3–4/7 PASS**: May be appropriate for investigative leads but not primary evidence
- **5–6/7 PASS**: May be appropriate for operational use with documented limitations
- **7/7 PASS**: Meets minimal evidence threshold for full forensic deployment

Important caveat: passing all seven requirements establishes *minimal* evidence for responsible deployment, not a guarantee of perfection. Ongoing monitoring, incident reporting, and periodic revalidation remain essential.

*2) Failure Mode Handling:* If a system fails specific requirements, targeted remediation may be possible:

- **Fail Meaningful Information / Accuracy**: Try different XAI method (e.g., switch from SHAP to Integrated Gradients)
- **Fail Error Rates**: Conduct stratified analysis to identify conditions where performance degrades
- **Fail Standards**: Establish pre-registered protocol for future validation studies
- **Fail Oversight**: Develop calibrated confidence scores using conformal prediction or Bayesian methods

The template provides a systematic pathway for identifying deficiencies and targeting improvements—transforming compliance from a binary pass/fail into a continuous quality improvement process.

## VI. DISCUSSION AND POLICY IMPLICATIONS

The analysis reveals a fundamental mismatch between regulatory intent and technical practice. Legal frameworks mandate explainability for high-stakes AI systems, yet current XAI deployment lacks the validation foundations that regulators, courts, and practitioners require. This section discusses implications for key stakeholders and proposes actionable recommendations.

### A. For Regulators and Standards Bodies

*1) Gap Identified:* Regulatory language—GDPR's "meaningful information," the AI Act's "appropriate accuracy"—lacks technical operationalization. This ambiguity enables "checkbox compliance" where systems generate explanations without validating their quality.

*2) Recommendations:* **Establish Technical Standards**. Regulatory bodies (e.g., EU AI Office, NIST) should publish technical standards specifying minimal evidence requirements for XAI validation, analogous to existing standards for DNA analysis or digital forensics. The seven requirements and thresholds proposed here provide a starting point.

**Mandate Validation Protocols**. Require pre-registered validation protocols with published benchmarks before high-risk AI systems can be deployed. Protocols should specify acceptance thresholds before data collection to prevent post-hoc p-hacking.

**Require Error Rate Disclosure**. Mandate that deployed systems document known failure modes and conditional error rates (e.g., explanation faithfulness by demographic group, image quality, score range). This mirrors Daubert's error rate requirement and enables risk-informed deployment.

**Periodic Revalidation**. Establish timelines for revalidation (e.g., annually) as models, methods, and datasets evolve. Face recognition systems aren't static—validation cannot be one-time certification.

**Demographic Fairness Requirements**. Extend validation requirements to include fairness thresholds—explanation faithfulness must meet minimal thresholds for all demographic groups, not just aggregate populations. As illustrated in our Grad-CAM example (Section 5.2), faithfulness for dark-skinned females falls to $\rho = 0.64$ compared to $\rho = 0.68$ overall—a disparity that exacerbates existing bias concerns in face recognition systems [2].

**Precedent**: The European Union's Medical Device Regulation (MDR 2017/745) [19] provides a model for risk-based AI oversight with technical standards, conformity assessment, and post-market surveillance. Adapting these principles to AI explainability could establish rigorous governance.

### B. For System Developers and Vendors

*1) Gap Identified:* Current development practices treat explanation generation as a feature add-on, not a core system requirement with validation obligations.

*2) Recommendations:* **Validation-First Development**. Incorporate faithfulness validation into the development lifecycle from the start, not as a post-deployment afterthought. XAI methods should be selected based on empirical validation performance, not popularity or visual appeal.

**Benchmark Participation**. Contribute to community development of standardized XAI benchmarks with ground truth. Publish validation results to establish credibility and enable comparative evaluation.

**Uncertainty Quantification**. Provide calibrated confidence intervals for explanations using conformal prediction [18] or Bayesian methods. Operators need uncertainty estimates to calibrate trust appropriately.

**Per-Instance Quality Scores**. Develop and deploy reliability indicators that enable operators to identify when specific explanations are unreliable. Aggregate validation metrics are insufficient—operators need case-level guidance. The Grad-CAM example's AUC of 0.71 shows this remains challenging but achievable.

**Transparent Limitation Documentation**. Clearly communicate known failure modes in user documentation and system interfaces. When an image falls into a documented failure mode (e.g., profile face, low resolution), flag this for operators automatically.

**Open Validation**. Publish validation protocols and results in peer-reviewed venues. Proprietary systems can be validated on public benchmarks without disclosing model weights.

**Business Case**: While validation adds development costs, it mitigates liability risks. Systems that contribute to wrongful arrests or fail Daubert challenges expose vendors to lawsuits. Proactive validation provides defensible due diligence.

### C. For Auditors and Oversight Bodies

*1) Gap Identified:* Auditors tasked with assessing AI system compliance lack technical tools and standards to evaluate explanation quality.

*2) Recommendations:* **Adopt Standardized Evaluation Protocols**. Use the compliance template (Section 5) or similar structured frameworks to systematically assess XAI validation. Require vendors to provide completed templates as part of procurement or compliance review.

**Independent Validation**. Don't rely solely on vendor-provided validation studies. Conduct independent testing on held-out datasets, particularly for high-stakes deployments. The difference between vendor claims and independent assessment can be substantial.

**Red Team Testing**. Employ adversarial evaluation to identify conditions under which explanations fail. Test edge cases: demographic groups underrepresented in training data, challenging poses, adversarial perturbations. The five failure modes identified in the Grad-CAM example emerged precisely from this kind of stratified analysis.

**Ongoing Monitoring**. Compliance isn't binary or static. Establish continuous monitoring programs that track explanation quality metrics over time as systems evolve and operational conditions change.

**Transparency Requirements**. Require that systems undergoing audit provide sufficient access for replication—validation datasets, model APIs (even if weights remain proprietary), and detailed methodology documentation.

**Precedent**: Financial services auditing (e.g., SOX compliance for algorithmic trading) provides models for independent technical evaluation of complex systems with legal accountability.

### D. For Courts and Legal Professionals

*1) Gap Identified:* Judges and attorneys lack technical expertise to evaluate XAI evidence presented in criminal proceedings, leading to either uncritical acceptance or blanket exclusion.

*2) Recommendations:* **Daubert Challenges for XAI Evidence**. When face recognition explanations are introduced as evidence, defense attorneys should challenge admissibility under Daubert. Critical questions include: Has the XAI method been validated with documented error rates across different conditions? Are published standards controlling its operation, or does deployment rely on ad-hoc vendor practices? Can the explanation be tested through controlled experiments that could potentially refute its claims? Whether peer review has occurred matters—but more importantly, whether that review addressed validation rather than just method description.

The Grad-CAM example's failure on Standards (no pre-registration), Meaningful Information ($\rho = 0.68$), and Accuracy (76%) would provide grounds for challenge.

**Expert Witness Standards**. Courts should require that expert witnesses presenting XAI evidence have conducted (or reviewed) rigorous validation studies, not merely familiarity with the XAI tool. An expert testifying "we used Grad-CAM" without validation data should face cross-examination on faithfulness, error rates, and failure modes.

**Judicial Education**. Provide training for judges on XAI fundamentals and validation principles through judicial education programs (e.g., Federal Judicial Center). Enable informed gatekeeping without requiring deep technical expertise. Simple questions—"What is the correlation between predicted and actual score changes?" "What percentage of cases fall into known failure modes?"—can reveal whether proper validation occurred.

**Standard Jury Instructions**. Develop model jury instructions for cases involving face recognition evidence. Jurors need to understand the critical distinction between model accuracy (how often the system correctly matches faces) and explanation accuracy (whether the highlighted features actually drove the decision). Validation metrics like correlation coefficients and error rates should be explained in accessible terms. Known limitations—such as degraded performance on profile faces or low-resolution images—deserve explicit mention, as does guidance on the appropriate weight XAI evidence should carry relative to other forensic methods.

**Precedent Development**. As cases involving XAI evidence accumulate, appellate decisions should establish precedent on admissibility standards, clarifying how Daubert applies to explainability methods specifically.

Key point: technical faithfulness is necessary but not sufficient for legal admissibility. Even validated explanations must be probative, reliable in the specific case context, and not unduly prejudicial.

### E. Who Benefits From Validated XAI?

The proposed validation framework serves multiple stakeholders with aligned interests in accuracy and accountability:

**Defendants and Accused Persons**. Validated explanations enable effective challenges to face recognition evidence. If an explanation fails validation thresholds, defense attorneys have grounds to argue for exclusion or reduced evidentiary weight.

**Law Enforcement and Forensic Analysts**. Validated systems protect agencies from liability risks associated with

wrongful arrests. Knowing when explanations are reliable versus unreliable enables more effective investigations and resource allocation. The 23% rejection rate in the Grad-CAM example means analysts can focus expert review where it's most needed.

**Regulatory Agencies**. Validated systems provide clear compliance evidence, reducing enforcement ambiguity and enabling risk-based oversight prioritization.

**System Developers**. Validation standards create level playing fields and enable differentiation based on empirical performance rather than marketing claims.

**Judges and Courts**. Validated evidence reduces Daubert hearing complexity and provides clear admissibility criteria, streamlining proceedings.

**Society**. Reduced wrongful identifications protect civil liberties; transparent accountability mechanisms build public trust in beneficial uses of face recognition technology.

### F. Remaining Gaps and Future Directions

While this article provides a framework for operationalizing existing regulatory requirements, several gaps require ongoing attention:

**Threshold Consensus**. The proposed thresholds ($\rho \geq 0.70$, 80% accuracy, etc.) are informed by statistical practice and analogous domains but require community consensus through standards development processes (ISO, NIST, professional societies).

**Dynamic Adaptation**. Validation standards must evolve as XAI methods, face verification architectures, and adversarial threats develop. Static standards risk obsolescence.

**Cross-Jurisdictional Harmonization**. U.S., EU, and other jurisdictions have different legal frameworks. International standards harmonization could reduce compliance complexity for multinational deployments.

**Fairness Integration**. Current regulatory frameworks address explainability and accuracy but lack explicit fairness requirements. Future standards should mandate that validation thresholds are met across demographic groups. Our Grad-CAM example revealed this gap concretely: while aggregate faithfulness reached $\rho = 0.68$, dark-skinned females experienced $\rho = 0.64$—falling further below the 0.70 threshold and compounding algorithmic bias concerns.

**Alternative Explanation Paradigms**. This article focuses on attribution-based XAI (saliency maps). Other paradigms—example-based explanations, concept-based interpretability, natural language rationales—require separate validation frameworks.

### G. A Call for Evidence-Based Policy

Current XAI practice has operated in a normative vacuum—researchers develop methods based on intuition, vendors deploy based on demand, and regulators mandate explainability without technical specificity. This article proposes a shift toward **evidence-based explainability policy**:

- Requirements grounded in measurable criteria
- Validation following scientific method principles

- Standards informed by empirical performance data
- Ongoing evaluation as systems and threats evolve

This mirrors the evolution of other forensic domains. DNA analysis, fingerprint comparison, and ballistic matching once lacked rigorous scientific foundations. Following high-profile wrongful convictions and critical reports (e.g., the 2009 NRC report on forensic science [14]), these fields developed validation protocols, error rate disclosure requirements, and proficiency testing standards.

Face recognition XAI stands at a similar inflection point. Documented wrongful arrests [3], [4] and regulatory mandates [8]–[10] create urgency for evidence-based standards. The framework proposed here—seven evidentiary requirements with operationalized thresholds and a compliance template—provides a starting point, not a final answer. Refinement through multi-stakeholder collaboration (researchers, practitioners, regulators, civil liberties advocates) is essential.

But the status quo—deploying explanations without validation—is scientifically indefensible and legally untenable. The time for evidence-based explainability policy is now.

## VII. CONCLUSION

Face recognition systems deployed in law enforcement operate at the intersection of impressive technical capabilities and profound accountability challenges. These systems achieve high matching accuracy—often exceeding 99.7% on benchmark datasets—yet their decision-making processes remain opaque. Explainable AI methods offer a path toward transparency by generating visual attributions highlighting influential facial features. However, current practice exhibits a critical gap: explanations are generated without rigorous validation of their faithfulness to model reasoning.

This gap matters because regulatory frameworks increasingly mandate not just explanations, but *accurate* explanations. The EU AI Act requires "accurate, accessible, and comprehensible information"; GDPR demands "meaningful information about the logic involved"; and U.S. courts applying Daubert standards require testable methods with known error rates. Current XAI practice—producing explanations without validating them—cannot demonstrate compliance with these requirements.

Through systematic analysis of three major regulatory frameworks (EU AI Act, GDPR, Daubert standard), we identified seven core evidentiary requirements: meaningful information, testability, known error rates, appropriate accuracy, adherence to standards, comprehensibility, and human oversight support. For each requirement, we proposed minimal technical evidence specifications, validation methods, and acceptance thresholds that operationalize vague legal concepts into measurable criteria.

The proposed framework reveals a troubling pattern: current practice achieves compliance in form but not substance. Systems generate explanations (satisfying literal regulatory language) without validation (failing the policy intent). This form-versus-substance gap exposes legal systems, defendants, and agencies to serious risks—wrongful identifications based on misleading explanations, Daubert inadmissibility challenges

that derail prosecutions, and regulatory enforcement uncertainty that creates legal exposure for deploying agencies.

The path forward requires evidence-based policy. We conclude with concrete recommendations:

**For Regulators**: Establish technical standards operationalizing vague legal language ("meaningful information," "appropriate accuracy") into measurable criteria. Mandate pre-registered validation protocols with published benchmarks. Require error rate disclosure including demographic stratification. Establish periodic revalidation requirements as systems evolve.

**For Developers**: Adopt validation-first development practices where XAI methods are selected based on empirical performance, not popularity. Contribute to community benchmark development. Provide calibrated uncertainty estimates and per-instance quality scores. Transparently document limitations and known failure modes. Publish validation protocols and results in peer-reviewed venues.

**For Auditors**: Adopt standardized evaluation protocols like the compliance template proposed here. Conduct independent validation beyond vendor claims. Employ red team testing for edge cases and demographic subgroups. Establish continuous monitoring programs. Require transparency enabling replication.

**For Courts**: Subject XAI evidence to rigorous Daubert scrutiny—testability, error rates, published standards, peer review. Require expert witnesses to demonstrate validation, not merely familiarity with tools. Develop judicial education programs on XAI validation principles. Create standard jury instructions explaining the distinction between model accuracy and explanation accuracy.

The compliance template and example validation (Grad-CAM on ArcFace) demonstrate both the feasibility and necessity of systematic assessment. The example system passed 3/7 requirements—sufficient for investigative leads under supervision, but not for primary evidence in legal proceedings. This nuanced assessment, grounded in measurable criteria, enables risk-informed deployment decisions that balance innovation with accountability.

Face recognition XAI stands where DNA analysis stood decades ago—at an inflection point between ad-hoc practice and scientific rigor. DNA analysis evolved from a novel forensic tool with uncertain reliability into a cornerstone of criminal justice, but only after developing validation protocols, error rate disclosure requirements, and proficiency testing standards. This evolution followed high-profile wrongful convictions and critical National Research Council reports demanding scientific foundations for forensic methods.

Face recognition XAI faces similar pressures. Robert Williams and Michael Oliver—arrested based on false matches—represent visible failures of a broader accountability gap. The EU AI Act, GDPR Article 22, and Daubert standards create legal obligations. Yet between regulatory mandate and technical reality lies a translation problem: how to operationalize legal requirements into measurable technical criteria?

This article provides that translation. The seven evidentiary requirements, minimal thresholds, and compliance template transform abstract legal concepts into concrete technical specifications. These specifications aren't final answers—threshold values require community consensus through standards development processes, and validation methods will evolve as XAI techniques advance. But they provide a starting point grounded in statistical practice, forensic science precedent, and existing regulatory frameworks.

The status quo—deploying explanations without validation—is scientifically indefensible and legally untenable. Explanations that are 68% faithful (below our 0.70 threshold) and 76% accurate (below our 80% threshold) may be better than nothing, but they're insufficient for contexts where liberty is at stake. The documented failure modes—profile faces, low resolution, occlusion, demographic disparities, borderline scores—reveal systematic patterns that operators must understand to make informed decisions.

The framework proposed here serves multiple stakeholders with aligned interests in accuracy and accountability. Defendants gain tools to challenge unreliable evidence. Law enforcement agencies protect themselves from liability while improving investigative effectiveness. Regulators obtain clear compliance criteria. Developers establish level playing fields. Courts streamline admissibility determinations. Society benefits from reduced wrongful identifications and increased trust in legitimate applications of face recognition technology.

Achieving these benefits requires multi-stakeholder collaboration—which is both urgent and achievable. Standards bodies like ISO and NIST must convene researchers, practitioners, regulators, and civil liberties advocates to refine thresholds and develop consensus benchmarks. We observe that research communities are beginning to shift from prioritizing subjective interpretability to objective faithfulness validation, though publication incentives still favor novel methods over rigorous evaluation. Vendors face perhaps the hardest challenge: embracing validation-first development even when it reveals uncomfortable limitations about their products. Courts, finally, need expertise to evaluate technical evidence rigorously—moving beyond uncritical deference to expert testimony toward informed gatekeeping.

This collaboration is both urgent and achievable. The technical foundations exist—counterfactual validation, conformal prediction, ground truth benchmarks. The legal frameworks create obligation and motivation. The wrongful arrest cases demonstrate tangible harms of inaction. What remains is translation and implementation.

This article provides the translation: seven requirements, minimal thresholds, a compliance template, and stakeholder-specific recommendations. Implementation requires commitment from regulators, developers, auditors, and courts to establish evidence-based validation as the norm rather than the exception.

The choice is clear. We can continue deploying explanations without validation, hoping that systems are trustworthy while lacking tools to verify that trust. Or we can demand evidence—testable predictions, known error rates, published standards, peer-reviewed protocols. The former preserves the status quo and its attendant risks. The latter builds accountability into AI systems from the foundation.

Face recognition technology offers genuine benefits for public safety and security. But those benefits cannot come at

the cost of civil liberties or scientific integrity. Validated explainability bridges that gap—enabling beneficial applications while providing the accountability mechanisms that protect individual rights. The framework proposed here shows the path. Now comes the hard work of walking it.

TABLES

TABLE I
REGULATORY REQUIREMENTS VS. CURRENT XAI PRACTICE

| Requirement | Current Practice | Gap | Im |
|---|---|---|---|
| **Meaningful Information** (GDPR Art. 22) | Visual saliency maps (Grad-CAM, SHAP) produced without validation | No verification that highlighted regions actually influenced decision | Inc tes tio tio |
| **Appropriate Transparency & Accuracy** (AI Act Art. 13) | Documentation describes XAI method used | No evidence that explanations are *accurate* representations of model reasoning | Op abl rec |
| **Testability** (Daubert) | XAI methods produce outputs but lack falsifiable hypotheses | Explanations cannot be empirically tested or refuted through controlled experiments | Fai sta den |
| **Known Error Rates** (Daubert; AI Act Art. 14) | Error rates reported for face verification accuracy, not for explanation faithfulness | No quantified failure modes of attribution methods; investigators don't know when explanations are unreliable | Ca exp exp |
| **Standards** (Daubert) | Ad-hoc deployment of XAI methods without published protocols or acceptance thresholds | No consensus standards for when explanation quality is sufficient for forensic use | Inc cie par |
| **Appropriate Accuracy** (AI Act Art. 13) | Verification models report accuracy metrics, but explanation accuracy is assumed, not measured | Attribution methods may systematically misidentify important features (studies show 40–69% accuracy) | Hi gua fal cat |
| **Human Oversight** (AI Act Art. 14) | Operators review XAI outputs without tools to assess explanation quality | Operators lack meta-information about explanation reliability for specific cases | Ca unn bec sta |

TABLE II
MINIMAL EVIDENCE REQUIREMENTS FOR XAI COMPLIANCE

| Requirement | Minimal Technical Evidence | Validation Method | Acceptance Threshold | Reporting Format |
|---|---|---|---|---|
| **Meaningful Information** (GDPR Art. 22) | Faithful attribution map where highlighted regions actually influence model decision | Counterfactual score prediction: $\Delta s_{\text{predicted}}$ vs. $\Delta s_{\text{actual}}$ | Pearson $\rho \geq 0.70$ between predicted and actual score changes | "Attribution faithfulness: $\rho = X.XX$ [95% CI: X.XX–X.XX]" |
| **Testability** (Daubert) | Falsifiable hypothesis about feature importance that can be empirically tested | Perturbation experiments with statistical hypothesis testing | $p < 0.05$ for $H_0$: attribution is random guessing; Cohen's $d \geq 0.5$ | "Testability: $\chi^2 = XX$, $p < 0.001$; attributions significantly predict score changes" |
| **Known Error Rates** (Daubert; AI Act Art. 14) | (1) Confidence interval calibration for predictions; (2) Documented failure modes | Conformal prediction for CI coverage; stratified evaluation by demographics/conditions | (1) 90–95% coverage for stated CIs; (2) Complete failure mode documentation | "CI calibration: 92% coverage at 90% CI. Known failure modes: [list]. Rejection rate: X%" |
| **Appropriate Accuracy** (AI Act Art. 13) | Quantified explanation accuracy independent of model accuracy | Ground truth test cases where true feature importance is known | Explanation accuracy $\geq$ 80% on ground truth benchmarks | "Explanation accuracy: 85% correct feature identification [benchmark: controlled perturbation suite]" |
| **Standards** (Daubert) | Pre-registered validation protocol with published acceptance criteria | Peer-reviewed validation study using standardized benchmark | Methods published in peer-reviewed venue; benchmark publicly available | "Validation protocol: [citation]. Benchmark: [name]. Results: [metrics]" |
| **Comprehensibility** (AI Act Art. 13) | Explanation + uncertainty quantification + limitations documentation | User study or expert evaluation of comprehensibility (secondary to technical faithfulness) | Target audience can correctly interpret explanation's meaning and limitations $\geq$ 75% of time | "Comprehensibility: XX% correct interpretation by [target audience] in controlled study" |
| **Human Oversight** (AI Act Art. 14) | Meta-level reliability indicator for each explanation (per-instance quality score) | Prediction confidence calibrated to actual accuracy on held-out validation set | Operator can discriminate between reliable/unreliable explanations with AUC $\geq 0.75$ | "Reliability indicator: AUC = 0.XX for predicting explanation error" |

## REFERENCES

[1] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "ArcFace: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4690–4699.

[2] P. Grother, M. Ngan, and K. Hanaoka, "Face recognition vendor test (FRVT) part 3: Demographic effects," National Institute of Standards and Technology (NIST), Gaithersburg, MD, Tech. Rep. NIST IR 8280, 2019.

[3] K. Hill, "Wrongfully accused by an algorithm," The New York Times, June 2020, retrieved from https://www.nytimes.com.

[4] ——, "Another arrest, and jail time, due to a bad facial recognition match," The New York Times, August 2023, retrieved from https://www.nytimes.com.

[5] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 618–626.

[6] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 4765–4774.

[7] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim, "Sanity checks for saliency maps," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2018, pp. 9505–9515.

[8] European Union, "Regulation (EU) 2024/1689 of the European Parliament and of the Council on Artificial Intelligence (AI Act)," Official Journal of the European Union, 2024, oJ L, 2024/1689.

[9] ——, "Regulation (EU) 2016/679 of the European Parliament and of the Council (General Data Protection Regulation)," Official Journal of the European Union, 2016, oJ L 119, 4.5.2016.

[10] Supreme Court of the United States, "Daubert v. Merrell Dow Pharmaceuticals, Inc." 509 U.S. 579, 1993, case No. 92-102.

[11] M. E. Kaminski, "The right to explanation, explained," *Berkeley Technology Law Journal*, vol. 34, no. 1, pp. 189–218, 2019.

[12] P.-J. Kindermans, S. Hooker, J. Adebayo, M. Alber, K. T. Schütt, S. Dähne, D. Erhan, and B. Kim, "The (un)reliability of saliency methods," in *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*.   Springer, 2019, pp. 267–280.

[13] Federal Rules of Evidence, "Federal Rules of Evidence, Rule 702: Testimony by Expert Witnesses," 28 U.S.C., 2011.

[14] National Research Council, *Strengthening Forensic Science in the United States: A Path Forward*.   Washington, DC: National Academies Press, 2009.

[15] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2017, pp. 3319–3328.

[16] S. Wachter, B. Mittelstadt, and C. Russell, "Counterfactual explanations without opening the black box: Automated decisions and the GDPR," *Harvard Journal of Law & Technology*, vol. 31, no. 2, pp. 841–887, 2017.

[17] J. Cohen, "Statistical power analysis for the behavioral sciences," *Lawrence Erlbaum Associates*, 1988.

[18] V. Vovk, A. Gammerman, and G. Shafer, "Algorithmic learning in a random world," 2005.

[19] European Union, "Regulation (EU) 2017/745 of the European Parliament and of the Council on Medical Devices," Official Journal of the European Union, 2017, oJ L 117, 5.5.2017.