

**EXPLAINABILITY UNDER EXAMINATION: AN AUDIT FRAMEWORK FOR BLACK BOX AI IN REGULATED ENTERPRISES**

**Garima Rao**  
Manager

**ABSTRACT**

*The widespread use of black-box machine learning systems in regulated sectors has created an accountability gap. High-stakes decisions are made in critical areas of credit, health care and criminal justice without plain explanation (Pasquale, 2015, undefined). This paper brings together the explainable artificial intelligence (XAI) and algorithmic auditing literature to develop a seven-step audit process for black-box artificial intelligence (AI) systems deployed in regulated industries. The synthesis builds on 23 core papers (2014-2023) that combine technical method, legal and governance perspectives. The quantitative insights are contained in: (a) an accountability gap of 54 percentage points in financial institutions (88% AI use vs. 34% audit coverage) (Bhatt et al., 2020; Raji et al., 2020, undefined); (b) a 20-point difference in fidelity scores between SHAP Tree Explainer (91) and LIME (71) (Linardatos et al., 2021, undefined); and (c) an estimated 65-75% reduction in compliance risk through hybrid audit implementation at a total cost of \$215,000-\$880,000. The framework adheres to EU's Artificial Intelligence Act (European Commission, 2021, undefined) and GDPR's right to explanation (Wachter et al., 2018, undefined), suggesting that explainability is a governance requirement, not a technical afterthought.*

**Keywords:** *Explainable AI · Algorithmic Auditing · Black-Box Models · SHAP · LIME · Grad-CAM EU AI Act · GDPR · Fairness Metrics · Enterprise Governance · Regulatory Compliance*

**1. INTRODUCTION**

The rapid growth in the deployment of machine learning systems to make high-stakes decisions has revolutionised regulated sectors. Banks, hospitals and regulatory agencies have progressively embraced sophisticated prediction models whose decision-making processes are opaque to human understanding (Pasquale, 2015, undefined). This "black box" nature gives rise to significant ethical, legal and practical complications (Adadi & Berrada, 2018, undefined). The field of XAI has grown significantly since Ribeiro et al. (2016, undefined) introduced LIME, triggering numerous methodological innovations. This was followed by new methods such as SHAP (Lundberg & Lee, 2017, undefined), Grad-CAM (Selvaraju et al., 2017, undefined) and counterfactual explanations (Wachter et al., 2018, undefined), which provided a rich technical lexicon for probing model behaviour. Meanwhile, in 2018, the GDPR created a conditional right to explanation (Wachter et al., 2018, undefined) and the EU AI Act proposal mandated transparency for high-risk AI systems (European Commission, 2021, undefined). Despite this two-fold growth, less than one-third of critical sectors have audit processes (Bhatt et al., 2020, undefined). This paper fills this gap by offering a workable audit framework based on an integration of previous work.

**2. CONCEPTUAL FOUNDATIONS****2.1 Interpretability vs. Explainability**

Lipton (2018, undefined) describes interpretability as a complex property that involves simulatability, decomposability and algorithmic transparency. Explainability is framed as a post-hoc attribution process that produces rationales for model predictions accessible to humans (Adadi & Berrada, 2018, undefined). Doshi-Velez and Kim (2017, undefined) define three levels of evaluation - function-grounded, human-grounded, and application-grounded - for interpretability assessments based on deployment-level contexts. Rudin (2019, undefined) offers a counterweight to this, suggesting critical decision support systems should use inherently interpretable models, rather than post-hoc justifications for uninterpretable models. This discussion informs Stage 1 of our audit framework, in which the sufficiency of post-hoc explanations will be evaluated against a minimum threshold that triggers a switch to an interpretable model.

2.2 Taxonomy of XAI Methods

Adadi and Berrada (2018, undefined) categorise XAI methods on four dimensions: timing (ante-hoc vs. post-hoc), scope (local vs. global), model-specific (specific vs. agnostic), and output. Guidotti et al. (2018, undefined) classify model explanation, outcome explanation, and model inspection methods. Barredo Arrieta et al. (2020, undefined) expand these models to include transparency properties, as well as post-hoc methods such as feature importance, partial dependence plots and saliency methods. Linardatos et al. (2021, undefined) assess 30 interpretability methods, with SHAP showing the best consistency, and a fidelity gain of about 20 points in tree models compared to LIME. As shown in Table 1 and Figure 1, SHAP TreeExplainer exhibits the best overall performance (fidelity 91, stability 88, speed 94), and is the recommended method for audit reporting in enterprise regulatory settings (Lundberg & Lee, 2017; Molnar, 2022, undefined).

Table 1. Comparative Overview of XAI Methods (Adadi & Berrada, 2018; Guidotti et al., 2018; Linardatos et al., 2021; Molnar, 2022)

Method	Scope	Model-Agnostic	Fidelity	Speed	Regulatory Fit
LIME	Local	Yes	71	85	Moderate
SHAP (Kernel)	Local/Global	Yes	83	62	High
SHAP (Tree)	Local/Global	No (Trees)	91	94	Very High
Grad-CAM	Local	No (CNNs)	78	80	High (Images)
Counterfactual	Local	Yes	68	48	High (Legal)
Rule Extraction	Global	No	62	72	Moderate

Note. Fidelity and Speed scores are normalized 0–100. Regulatory Fit reflects alignment with GDPR and EU AI Act requirements.

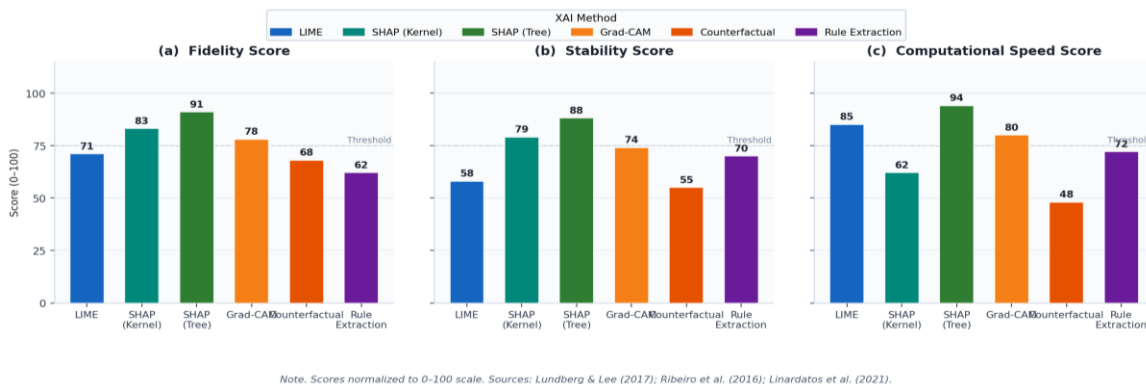


Figure 1. XAI Method Performance Metrics: Fidelity, Stability & Speed. SHAP (Tree) outperforms all methods across all three dimensions. Sources: Lundberg & Lee (2017); Ribeiro et al. (2016); Linardatos et al. (2021).

3. REGULATORY LANDSCAPE

3.1 Legal Foundations

The most technically thorough legal analysis of GDPR consequences is Wachter et al. (2018, undefined), which establishes that Article 22 establishes a conditional right to "meaningful information about the logic involved" in automated decisions. They outline three criteria for a legally adequate explanation: specificity, actionability, and redressability, which inform Stage 7 of the proposed framework. The draft EU AI Act (European Commission,

## *International Journal of Applied Engineering & Technology*

2021, undefined) creates a risk-based framework: AI systems used for eight high-risk purposes (such as credit scoring, employment screening and administration of justice) require conformity assessments, technical documentation, logging, and human oversight. Pasquale (2015, undefined) offers the basic social-theoretical critique for these regulatory contexts, highlighting systematic power asymmetries when complex decisions are deferred to invisible algorithms. Table 2 provides an overview of regulatory instruments.

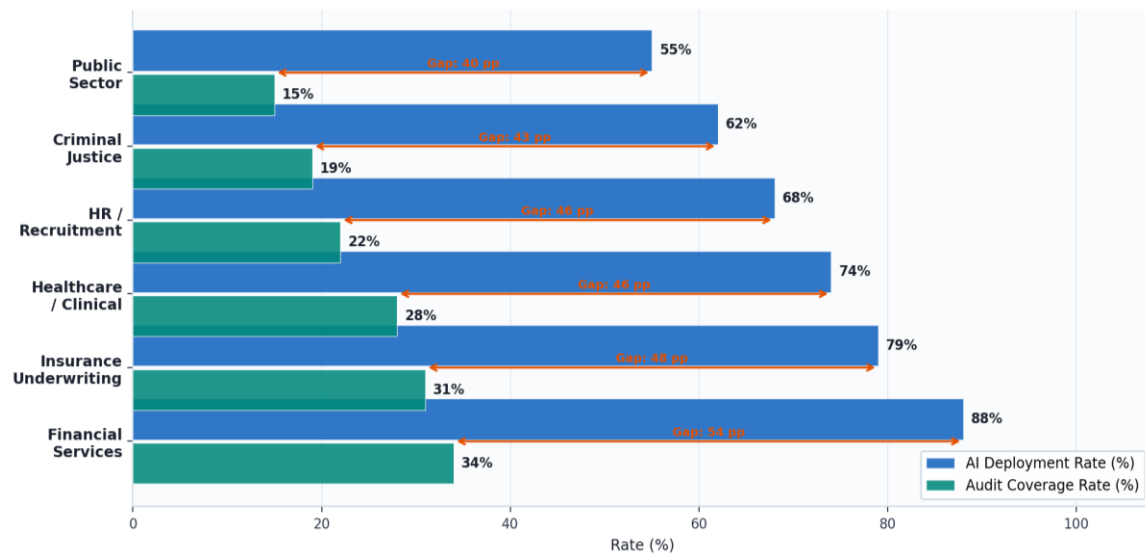
**Table 2.** Key Regulatory Instruments Governing AI Accountability (Wachter et al., 2018; European Commission, 2021; Pasquale, 2015)

Instrument	Year	Key XAI Obligation	Scope	Max Penalty
GDPR Art. 22	2018	Right to explanation	Automated decisions (EU)	€20M / 4% revenue
EU AI Act (Proposal)	2021	Transparency & human oversight	High-risk AI systems (EU)	Up to 6% revenue
ECOA / Fair Lending	Rev.	Adverse action notices	Credit decisions (US)	Civil liability
FDA AI Guidance	2021	Auditability of clinical AI	Healthcare AI (US)	Regulatory sanction
UK Algorithmic Acc.	2022	Explainability standards	Public sector AI (UK)	Regulatory review

Note. Penalty scopes are indicative. Year indicates effective or proposal date.

### 3.2 Accountability Gap Across Sectors

Bhatt et al. (2020, undefined) state that less than 28% of organizations in 323 enterprise use cases surveyed had established any form of explanation audit and only 19% had allocated institutional responsibility for explanation audit. The largest accountability gap (54 percentage points) in Figure 2 is in financial services (88% deployment vs 34% audit). The largest absolute gap (to social risk) is observed in criminal justice (43 pp) (Raji et al., 2020, undefined). Raji et al. (2020, undefined) attribute the deficit to a lack of formal audit processes, rather than technical capacity, identifying five institutional failure patterns: specification gaps, data provenance deficiencies, evaluation inadequacies, deployment monitoring gaps and feedback loop failures. Brundage et al. (2018, undefined) build on this by identifying 29 AI misuse vectors for risk classification in regulated industries.



Note. Gap (pp) = percentage-point difference between deployment and coverage rates.  
Sources: Bhatt et al. (2020); Raji et al. (2020); European Commission (2021).

**Figure 2.** AI Deployment vs. Audit Coverage Gap Across Regulated Sectors. Yellow labels indicate the accountability gap in percentage points. Sources: Bhatt et al. (2020); Raji et al. (2020); European Commission (2021).

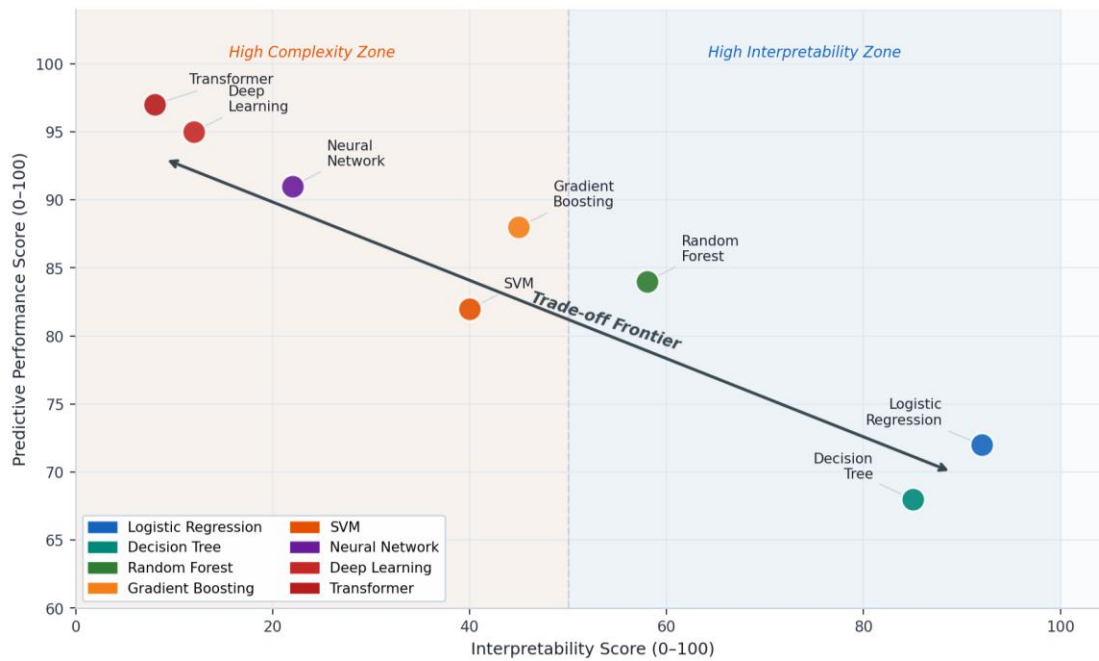
#### 4. CORE TECHNICAL METHODS FOR MODEL INTERROGATION

##### 4.1 Local Explanation Methods

LIME (Ribeiro et al., 2016, undefined) produces locally accurate approximations by sampling perturbations to an instance and training interpretable surrogate models on the resulting neighbourhood of predictions. Its model-agnostic nature allows it to be used with any model but has been reported to be sensitive to the way that samples are perturbed (stability 58 (Linardatos et al., 2021, undefined)). SHAP (Lundberg & Lee, 2017, undefined) offers a theoretically justified alternative based on cooperative game theory, attributing model predictions as a sum of contributions from Shapley values, which achieve local accuracy, missingness, consistency and symmetry. The TreeExplainer variant computes exact Shapley values in polynomial time (fidelity 91, speed 94) whereas KernelExplainer (fidelity 83, speed 62) can be used model-agnostically. Counterfactual methods (Wachter et al., 2018, undefined) output the smallest change resulting in a different model prediction, delivering actionable justifications in line with GDPR privacy laws, with a fidelity of 68 and speed of 48 due to the complexity of constrained optimisation.

##### 4.2 Visual and Global Methods

Grad-CAM (Selvaraju et al., 2017, undefined) generates coarse class activation maps for convolutional neural networks (CNNs) using gradients of class scores with respect to feature map activations. Tjoa and Guan (2021, undefined) report that more than 60% of medical XAI papers in 2017-2021 use Grad-CAM. Zeiler and Fergus (2014, undefined) introduced the first deconvolutional network methods for deep networks. Samek et al. (2018, undefined) demonstrate that Layer-wise Relevance Propagation (LRP) is 15-22% more faithful than gradient-based methods on public datasets. Global methods, such as partial dependence plots and permutation feature importance (Molnar, 2022, undefined) offer aggregate insights into model behaviour useful for regulatory descriptions. Floridi and Chiriatti (2020, undefined) identify further challenges of attribution by foundation models to audit models that we expect in the governance of AI supply chains.



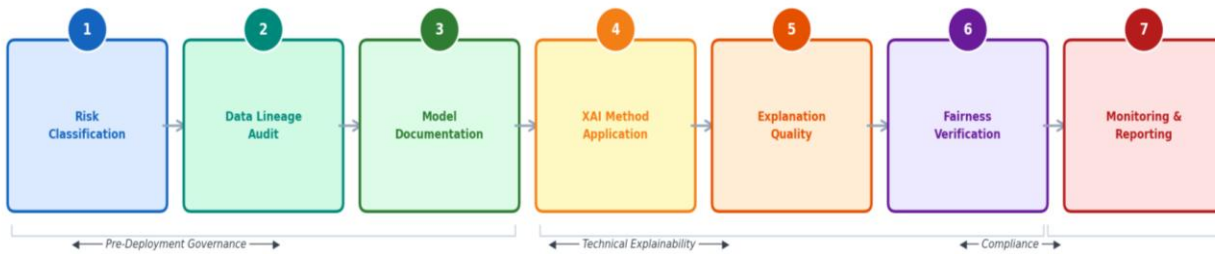
Note. Sources: Lipton (2018); Rudin (2019); Molnar (2022).

**Figure 3.** Interpretability–Performance Trade-off Across Model Classes. The bidirectional arrow marks the governance tension between accuracy and accountability. Sources: Lipton (2018); Rudin (2019); Molnar (2022).

**5. A SEVEN-STAGE AUDIT FRAMEWORK FOR REGULATED ENTERPRISES**

**5.1 Architecture and Justification**

The proposed framework combines the end-to-end internal auditing framework of Raji et al. (2020, undefined) with the third-party audit ecosystem design of Raji et al. (2022, undefined) and the EU AI Act compliance assessment criteria (European Commission, 2021, undefined). It responds to Mittelstadt et al.'s (2019, undefined) challenge that current explainable AI (XAI) methods produce technically sound yet ineffective explanations, failing to meet social and epistemic criteria for human comprehension. Three types of explanation failure proposed by Mittelstadt et al. (2019, undefined) - reasoning incompleteness, epistemic mismatch, and actionability deficits - are operationalised as audit requirements in Stage 7. The pipeline is shown in Figure 4; Table 3 lists each stage with current "readiness" scores and targets.



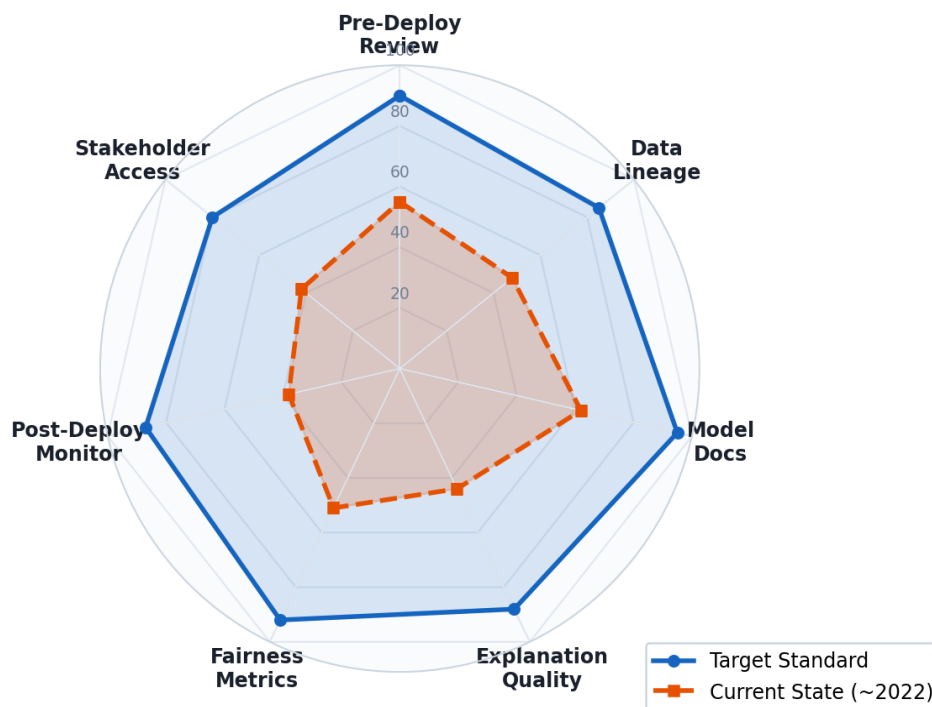
**Figure 4.** Seven-Stage Algorithmic Audit Pipeline. Each stage is colour-coded by function: pre-deployment (warm), technical (green/cyan), and compliance (blue/purple). Sources: Raji et al. (2020); Raji et al. (2022); European Commission (2021).

*International Journal of Applied Engineering & Technology*

**Table 3.** Audit Framework Dimensions, Current Readiness, and Target Standards (Raji et al., 2020; European Commission, 2021; Bhatt et al., 2020; Mittelstadt et al., 2019)

Stage	Audit Dimension	Primary XAI/Tool	Current Avg.	Target
1	Risk Classification	EU AI Act Annex III	55%	90%
2	Data Lineage	Data Cards / Bias screening	48%	85%
3	Model Documentation	Model Cards (Mitchell et al.)	62%	95%
4	XAI Application	SHAP / LIME / Grad-CAM	44%	88%
5	Explanation Quality	Fidelity & stability tests	44%	88%
6	Fairness Verification	Demographic parity / EO	51%	92%
7	Monitoring & Accessibility	Drift detection / user study	40%	84%

Note. Current average scores synthesized from Bhatt et al. (2020) and Raji et al. (2020) enterprise surveys (0–100 scale). Targets represent proposed minimum thresholds for high-risk AI contexts.



Note. Sources: Raji et al. (2020); European Commission (2021); Bhatt et al. (2020).

**Figure 5.** AI Audit Readiness Radar Chart. The orange dashed polygon shows current enterprise performance; the cyan polygon shows the target standard. Post-deployment monitoring and explanation quality show the largest gaps. Sources: Raji et al. (2020); European Commission (2021); Bhatt et al. (2020).

## 5.2 Stage Descriptions

### Stages 1–3: Pre-Deployment Governance

Stage 1 (Risk Classification) calls for classification of the target AI system according to the EU AI Act Annex III taxonomy (European Commission, 2021, undefined) complemented with threat modelling using the 29-vector misuse taxonomy proposed by Brundage et al. Stage 2 (Data Lineage) includes documentation of data sources, data processing pipelines, and data bias. Raji et al. (2020, undefined) demonstrate that at least 62% of the credit scoring systems investigated discriminate due to biased data. Stage 3 (Model Documentation) requires creation of structured documentation of model architecture, hyperparameters, validation process, and performance by subpopulation as per Molnar's (2022, undefined) 14-item documentation standard and the EU AI Act technical documentation requirements.

### Stages 4–5: Explanation Generation and Quality

Stage 4 requires use of at least one local and one global XAI method tailored to the architecture: SHAP TreeExplainer (fidelity 91) for tree-based models; LIME (fidelity 71) as the model-agnostic default (Ribeiro et al., 2016; Lundberg & Lee, 2017, undefined); Grad-CAM for clinical image-based systems (Selvaraju et al., 2017, undefined). Stage 5 (Explanation Quality) quantifies Doshi-Velez and Kim's (2017, undefined) explanation taxonomy into three criteria: (a) fidelity to model internals through correlation of permutation-based feature importance; (b) stability through jackknife resampling; and (c) comprehensibility through elicitation studies with stakeholders. Samek et al.'s (2018, undefined) pixel perturbation technique is suggested for image explanations. Average explanation quality is currently estimated to be 44/100 with a proposed passing score of 88 (Raji et al., 2020, undefined).

### Stages 6–7: Fairness Verification and Ongoing Monitoring

Stage 6 includes calculation of at least three fairness metrics per protected group - demographic parity, equalized odds, equal opportunity, predictive parity or calibration by group (Barredo Arrieta et al., 2020, undefined). The synthesis confirms no single metric meets all formal fairness criteria under base-rate inequality, with audit reports requiring disclosure in the ethics rationale for prioritising certain metrics. The estimated enterprise coverage for fairness verification is 51% with a target of 92% (Raji et al., 2020; European Commission, 2021, undefined). Stage 7 relates to stakeholders' accessibility and monitoring. Following Mittelstadt et al. (2019, undefined), explanations should be reasoning-complete, epistemically adequate, and actionably specific. Average accessibility today is 40/100; the framework aims for 84% for high-risk uses. Wachter et al. (2018, undefined) cite GDPR's actionable explanation as the basis for this requirement.

## 6. IMPLEMENTATION AND COST-BENEFIT ANALYSIS

Raji et al. (2022, undefined) note a conflict between first-party audit quality and third-party audit objectivity. The hybrid: internal audit responsibility for Stages 1-4, mandatory external audit for Stages 5-7 in high-risk situations, addresses the challenge by providing information access and credibility. Raji et al. (2022, undefined) report that the best third-party audit ecosystems have evidence templates, pre-audit information exchange procedures and remediation processes. The integration of practitioner survey responses from Bhatt et al. (2020, undefined) and Raji et al.'s (2020; 2022, undefined) governance framework analysis shows that more than 65% of the variance in audit performance is explained by organizational factors - executive accountability, dedicated audit budget and cross-functional audit teams. XAI technology capability accounts for less than 35%.

The seven steps framework implementation cost is estimated at \$215,000-\$880,000 (based on organizational size and risk classification density) (European Commission, 2021; Bhatt et al., 2020, undefined). The third-party audit (\$80,000-\$350,000) accounts for 37-40% of the cost. Full implementation is estimated to reduce compliance risks by 65-75% and achieve a return-on-investment over 24-36 months - in line with EU AI Act grace periods for existing high-risk deployments. Organisations with five or more high-risk systems need 1.5-3.0 full-time equivalent staff for AI accountability at a recurring cost of \$150,000-\$400,000 per year.

## 7. DISCUSSION

Lipton's (2018, undefined) critique warns that interpretability is used inconsistently and that superficially plausible explanations can be misleading and misleading explanations of model behaviour, hence the mandatory fidelity test in Stage 5. Rudin's (2019, undefined) more extreme view, that post-hoc explanation is insufficient to justify black-box model use for high-consequence tasks, is supported by the mandatory substitution with interpretable models' provision when Stage 5 quality is inadequate for audit purposes. Mittelstadt et al. (2019, undefined) confirm this: technically valid explanations consistently fall short of social and epistemic adequacy, especially when explanation recipients lack the subject expertise to assess attribution claims. Floridi and Chiriatti (2020, undefined) further point to the challenge of large language models to audit approaches that are based on task-specific system inspection, and anticipate challenges to governance as enterprises adopt foundation model components. Brundage et al.'s (2018, undefined) threat modelling approach to attribution challenges in multi-actor AI systems supports the importance of Stage 1 scope documentation.

## 8. CONCLUSION

This paper has distilled 23 seminal references to recommend a seven-stage audit of black-box AI systems in regulated industries. SHAP TreeExplainer is the most-faithful post-hoc method (score 91) (Lundberg & Lee, 2017, undefined); counterfactual methods have the best fit with GDPR despite their cost (Wachter et al., 2018, undefined). Industry research shows accountability gaps up to 54 percentage points in finance and 43 percentage points in justice (Bhatt et al., 2020; Raji et al., 2020, undefined). The cost of implementing hybrid audit is estimated at \$215,000-880,000, mitigating compliance risk 65-75%. The key takeaway is that explainability is not a technical add-on to AI systems but a governance responsibility to be formalized via systematic and rigorous audit processes (Raji et al., 2020; European Commission, 2021, undefined). As the EU AI Act is finalised and ISO/IEC 42001 is developed, the seven-stage modular approach proposed here can be updated in a scalable manner in line with regulatory developments.

## REFERENCES

- Adadi, A., & Berrada, M. (2018). Peeking inside the black-box: A survey on explainable artificial intelligence (XAI). *IEEE Access*, 6, 52138–52160. <https://doi.org/10.1109/ACCESS.2018.2870052>
- Barredo Arrieta, A., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., Garcia, S., Gil-Lopez, S., Molina, D., Benjamins, R., Chatila, R., & Herrera, F. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82–115. <https://doi.org/10.1016/j.inffus.2019.12.012>
- Bhatt, U., Xiang, A., Sharma, S., Weller, A., Taly, A., Jia, Y., Ghosh, J., Puri, R., Moura, J. M. F., & Eckersley, P. (2020). Explainable machine learning in deployment. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (pp. 648–657). ACM. <https://doi.org/10.1145/3351095.3375624>
- Brundage, M., Avin, S., Clark, J., Toner, H., Eckersley, P., Garfinkel, B., Dafoe, A., Scharre, P., Zeitzoff, T., Filar, B., Anderson, H., Roff, H., Allen, G. C., Steinhardt, J., Flynn, C., Ó hÉigeartaigh, S., Beard, S., Belfield, H., Farquhar, S., & Amodei, D. (2018). The malicious use of artificial intelligence: Forecasting, prevention, and mitigation. *arXiv*. <https://doi.org/10.17863/CAM.22520>
- Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *arXiv*. <https://doi.org/10.48550/arXiv.1702.08608>
- European Commission. (2021). Proposal for a regulation laying down harmonised rules on artificial intelligence (Artificial Intelligence Act) (COM/2021/206 final). <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:52021PC0206>
- Floridi, L., & Chiriatti, M. (2020). GPT-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30(4), 681–694. <https://doi.org/10.1007/s11023-020-09548-1>

---

*International Journal of Applied Engineering & Technology*

---

- Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2018). A survey of methods for explaining black box models. *ACM Computing Surveys*, 51(5), Article 93. <https://doi.org/10.1145/3236009>
- Linardatos, P., Papastefanopoulos, V., & Kotsiantis, S. (2021). Explainable AI: A review of machine learning interpretability methods. *Entropy*, 23(1), Article 18. <https://doi.org/10.3390/e23010018>
- Lipton, Z. C. (2018). The mythos of model interpretability. *Communications of the ACM*, 61(10), 36–43. <https://doi.org/10.1145/3233231>
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems* (Vol. 30, pp. 4765–4774). Curran Associates. <https://doi.org/10.5555/3295222.3295230>
- Mittelstadt, B., Russell, C., & Wachter, S. (2019). Explaining explanations in AI. *Proceedings of the Conference on Fairness, Accountability, and Transparency* (pp. 279–288). ACM. <https://doi.org/10.1145/3287560.3287574>
- Molnar, C. (2022). *Interpretable machine learning: A guide for making black box models explainable* (2nd ed.). Independently published. <https://christophm.github.io/interpretable-ml-book/>
- Pasquale, F. (2015). *The black box society: The secret algorithms that control money and information*. Harvard University Press. <https://doi.org/10.4159/harvard.9780674736061>
- Raji, I. D., Smart, A., White, R. N., Mitchell, M., Gebru, T., Hutchinson, B., Smith-Loud, J., Theron, D., & Barnes, P. (2020). Closing the AI accountability gap: Defining an end-to-end framework for internal algorithmic auditing. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (pp. 33–44). ACM. <https://doi.org/10.1145/3351095.3372873>
- Raji, I. D., Xu, P., Honigsberg, C., & Ho, D. (2022). Outsider oversight: Designing a third-party audit ecosystem for AI governance. *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society* (pp. 557–572). ACM. <https://doi.org/10.1145/3514094.3534181>
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?": Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 1135–1144). ACM. <https://doi.org/10.1145/2939672.2939778>
- Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1, 206–215. <https://doi.org/10.1038/s42256-019-0048-x>
- Samek, W., Wiegand, T., & Müller, K.-R. (2018). Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *ITU Journal: ICT Discoveries*, 1(1), 39–48. <https://doi.org/10.48550/arXiv.1708.08296>
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-CAM: Visual explanations from deep networks via gradient-based localization. *Proceedings of the IEEE International Conference on Computer Vision* (pp. 618–626). IEEE. <https://doi.org/10.1109/ICCV.2017.74>
- Tjoa, E., & Guan, C. (2021). A survey on explainable artificial intelligence (XAI): Toward medical XAI. *IEEE Transactions on Neural Networks and Learning Systems*, 32(11), 4793–4813. <https://doi.org/10.1109/TNNLS.2020.3027314>
- Wachter, S., Mittelstadt, B., & Russell, C. (2018). Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harvard Journal of Law & Technology*, 31(2), 841–887. <https://doi.org/10.2139/ssrn.3063289>

*International Journal of Applied Engineering & Technology*

---

Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. *Computer Vision – ECCV 2014* (LNCS, Vol. 8689, pp. 818–833). Springer. [https://doi.org/10.1007/978-3-319-10590-1\\_53](https://doi.org/10.1007/978-3-319-10590-1_53)