*Original Article*

# Bridging Data Science and Compliance   Intelligent AML Systems by Design: Domain-Aware Machine Learning

Pratik Chawande
Independent Research, Dallas , Tx

**Abstract:** *The rapid increase in global financial transactions has made the old, rules based, Anti-Money laundering (AML) systems ineffective and efficient. Such systems produce too many false positives, are expensive to operate, and fail to keep pace with more advanced and progressive typologies of crime. The paper suggests an innovative architecture of the next generation AML systems by integrating the powerful machine learning (ML) with deep compliance knowledge domain. We support a domain-sensitive strategy in which the development of the ML model, which includes feature engineering and natural language processing (NLP), pattern recognition and anomaly detection, is domain-directed, meaning that it is driven by the knowledge in the financial crime field, regulatory reporting, and operational specificity of investment banking. The paper examines the technical architectures that combine unsupervised learning to detect anomalies, supervised models to detect typology, and NLP to generate alert narrative automatically and analysis of unstructured data. More importantly, it deals with the extremely important issues of model explainability and adversarial robustness and the combination of both with the current compliance processes in order to make them regulatory-acceptable and effective. This paradigm will increase the precision of detection, minimize the workload on investigators, and bring a more intelligent and lively approach to financial crime defense through the combination of the technical profundity of data science with the fine-tuning of compliance logic.*

**Keywords:** *Anti-Money Laundering (AML), Domain-Aware Machine Learning, Compliance Automation, Financial Crime Detection, Explainable AI (XAI), Natural Language Processing (NLP), Anomaly Detection, Regulatory Technology (Regtech), Investment Banking, Model Governance.*

## 1. Introduction

International pressure to fight money laundering and terrorist financing is a multi-trillion-dollar problem to financial institutions. The regulatory requirements are such that regulatory frameworks such as the Bank Secrecy Act (BSA), the EU Anti-Money laundering Directives (AMLD), and Financial Action Task Forces (FATF) recommendations impose strict customer due diligence (CDD) and monitoring of transactions. The traditional compliance systems mainly utilize non-dynamically based rules (e.g., flag all cash transactions over 10,000 or flag movement of funds between several accounts). Although this is transparent and auditable, the approach has a limitation. It produces alarming rates of false-positives which can reach above 95 percent and causes compliance teams to be overloaded with low-value alerts resulting in a severe level of operational fatigue. This fatigue of alert results in the lack of critical signals during the noise. Moreover, criminal justice systems are reactive in nature and thus can be easily bypassed by criminals who constantly refine their own models to bypass any known thresholds and patterns.

The introduction of big data analytics and artificial intelligence is an opportunity to change something. Machine learning (ML), which is capable of learning non-linear and more complex patterns by analyzing high-dimensional data sets very large in scale, has the potential to transform AML in the rules-based to a risk-based, predictive paradigm. Nevertheless, the direct use of generic, off-the-shelf ML algorithms on the compliance domain has left behind suboptimal results and regulatory cynicism. The root cause of the failure frequently lies in the lack of a fundamental understanding: data scientists might not understand the money laundering typologies, regulatory expectations, and the realities of the alert adjudication, whereas compliance experts might not comprehend the capabilities and the constraints of the ML models.

The research hypothesis presented in this paper is that successful integration of data science and compliance needs to be designed by creating Intelligent AML Systems with a background of Domain-Aware Machine Learning. Domain-aware ML is not just the use of algorithms on the financial data; it is instead the procedure of putting compliance logic, subject matter knowledge and business context directly into each phase of the ML lifecycle. These involve problem formulation, feature engineering, model selection, training, validation and deployment. It goes beyond pure algorithmic performance to include explainability, auditability, and non-

disruptive integration into human-managed compliance processes, especially those in such complex contexts as investment banking with complex products, international flows and clients.

The following parts of this paper will consider the architectural elements of such intelligent systems, describe the innovations in NLP and pattern recognition to ensure compliance, discuss the key imperative of explainability and governance and analyze the cross-domain applications and the future direction of domain-sensitive AI in financial services compliance.

## 2. Main Body

### 2.1. The Urgent Need to Change the AML Paradigm

The weaknesses of the old AML systems are well-known but it is worth revisiting them and remind the world of the necessity of the innovation. A normal big international bank spends more than 1 billion dollars every year on the compliance of financial crimes and staffs thousands of investigators, however, less than 1 percent of hot spots activity reports (SARS) result in additional action by a law enforcement agency. Such appalling inefficiency is the result of a number of fundamental concerns.
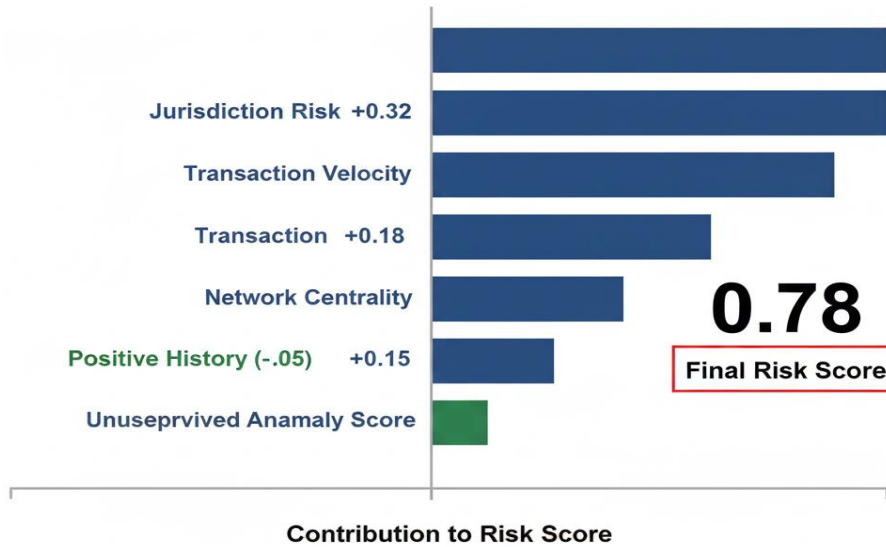


**Contribution to Risk Score**
**Fig 1: Shap Water Plot for Alert Explanation**

- Higher False Positives: Rules are simplistic and fail to put in consideration of valid contextual variations, where many normal behaviors of low-risk clients are flagged.
- Easily Hacked: Hackers research limits and patterns and organize the business to fly under radar (e.g., "smurfing") (Liu, 2025).
- Absence of Flexibility: Rule sets are hard to update, and they are not automated, which keeps the systems susceptible to new typologies over some time.
- Isolated Data: Rules tend to process individual data streams (e.g. transactions), without integrating signals of KYC data, news, network relationships and external watchlists.

Domain-aware ML suggests a solution to this by developing systems that learn a continuously varying contextual model of normal behavior per customer and relationship, are able to identify subtle deviations which predict illicit activity, and adapt as new threats emerge.

### 2.2. Domain-Aware Machine Learning Pillars of AML

An intelligent AML system design is based on three pillars that are interconnected and have compliance logic embedded in them.

#### 2.2.1. Domain-based feature engineering and Data synthesis

It is the most important adage of garbage in, garbage out. Among the raw transactional data (amount, date, counterparty), it is not enough. Domain-specific feature engineering designs predictive features that represent the concepts of compliance risk:

- Behavioral Characteristics: It is not a violation of a global standard but a deviation of a customer on his or her historical behavior (e.g., a sudden increase in the volume of transactions with a high-risk place).
- Network Features: Measures using transaction network graphs, including centrality, clustering coefficients, and the risk profile of the entities being connected (e.g. a client being connected with dozens of shell companies).
- Temporal Features: A transaction that is out of business hours is unusual (e.g., transactions are always out of business hours) and seasonality (Pérez, 2025).
- Product Risk Embeddings: Representing the underlying degree of riskiness of complex investment banking products (e.g., syndicated loans, derivatives, flows of prime brokerage) as model-readable features.

The close coordination between quants and AML analysts is required in this process such that features are not

only statistically significant but also meaningful and defensible in regulatory sense.

## 2.2.2. Hybrid Modeling Architecture

AML cannot be solved with the help of one algorithm. A sound system is one that has a hybrid, layered approach.
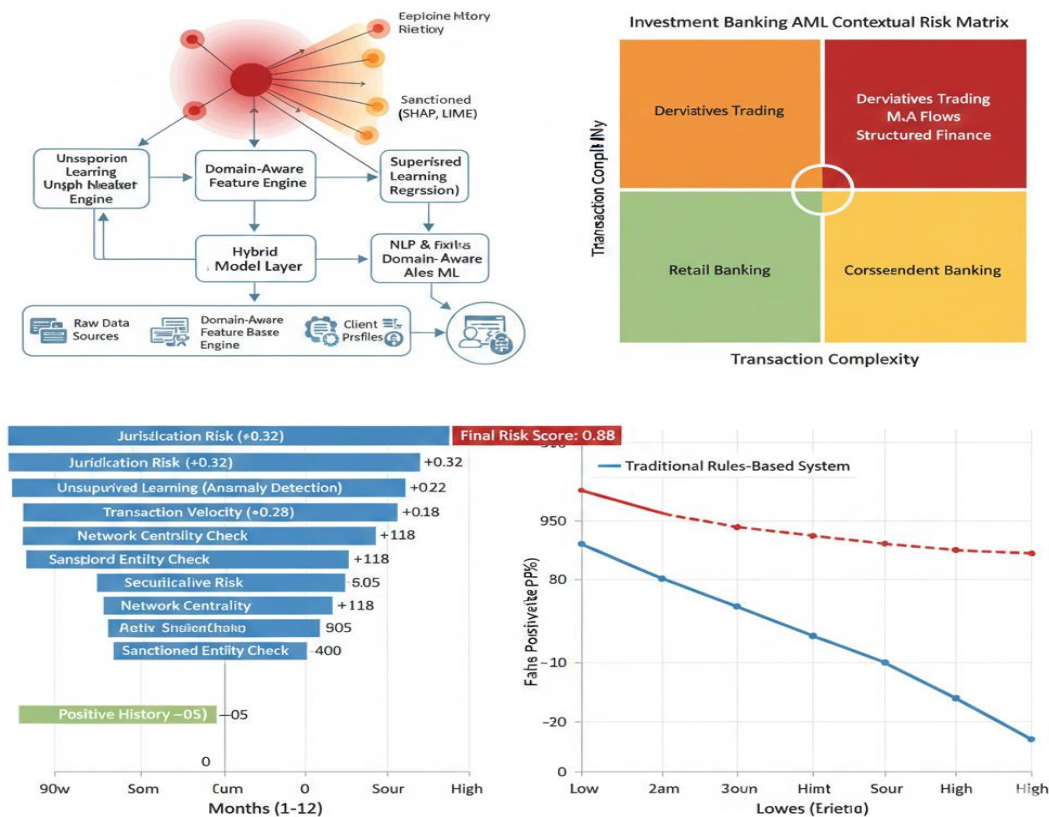


**Fig 2: Domain-Aware Hybrid AML System Architecture**

Human-in-the-Loop Alert Management Dashboard
- Unsupervised Learning (Anomaly Detection): Isolation Forests or Variational Autoencoders are models that are trained to learn a compact representation of normal activity, on a case-by-case basis, per client segment. They are good at identifying new and previously unknown suspicious patterns, which do not require labeled historical data, the unknown unknowns.
- Learned Supervision (Typology Classification): Alerts in the past that were eventually reported as SARS (true positives) and also those that were false positives are used to train a model such as XGBoost or Random Forests. They are taught to identify the complex feature combinations that are related to familiar laundering typologies (e.g., layering, structuring).
- Ensemble & Risk Scoring: Both the unsupervised and supervised models are combined to form a single risk score, which is more often than not weighted using domain-specific weights. This score puts alerts as a priority to investigators (Espinoza, 2025).

## 2.2.3. NLP to Automate and Enrich Compliance

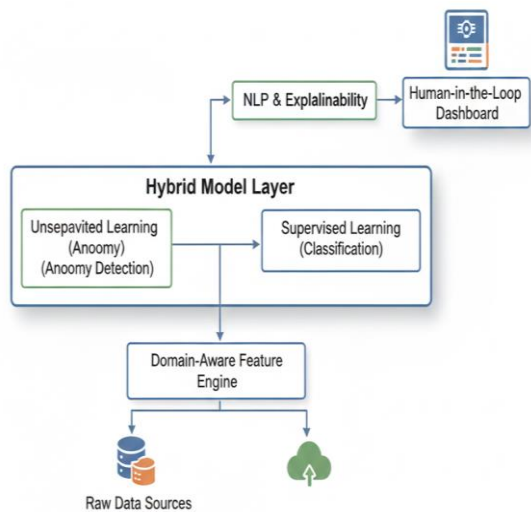NLP is the foundation of the domain-sensitive intelligence, that is, the intelligence going beyond numbers.
- Automated Alert Narrative Generation: As an alternative to giving an investigator a raw risk score and a list of transactions, a NLP module can produce a plain language summary: "Reminder received, because wire transfer volume to Jurisdiction Y increased 450 percent over the last 12 months, and because the client has received a negative news item about his parent company in the recent past. This saves a lot of time in review cases.
- Unstructured Data Analysis NLP models, especially Transformer-based models (such as BERT) trained on regulatory financial text will screen emails and chat messages sent by customers, as well as news articles, to detect risk factors (sentiment, reference to authorized entities, use of ambiguous language).
- KYC Profile Automation: Automated extraction and validation of entity data of corporate documents, beneficial ownership structures and biographical notes.

## 2.3. Patterning Recognition and Network Analytics

Criminality usually does not exist in individual behavior but in the relations between entities. This is formalised in domain-aware ML by graph machine learning.

- Dynamic Relationship Graphing: This involves building temporal graphs in which the nodes are clients/accounts and the edges are transactions, address shared, or director. Risk signals can be spread across the network by algorithms such as Graph Neural Networks (GNNs) in order to establish clusters of high-risk activity that would not otherwise be evident at the individual account level
- Temporal Pattern Mining: The specific sequence patterns identified to be related to the process of layering, (Hanif, 2025) i.e. quick cycles of deposits, transfers, and withdrawals on many accounts in a limited time frame.



**Fig 3: Risk Propagation in Network**

A (Simple diagram) would consist of a central low-risk node linked to a number of nodes. One of these secondary nodes is an indicator of high-risk. The propagation of a risk, through GNN, causes the risk to bleed back to a central node, causing a re-evaluation, which is an example of guilt-by-association in a financial network.)
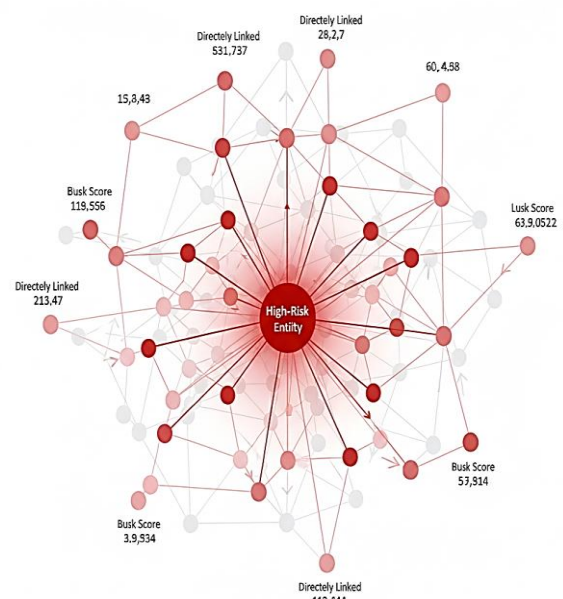
## 2.4. The Critical Bridge: Explainability, Adversarial Robustness and Model Governance

The main impediment to regulation and operation adoption is the black box problem. A domain-conscious system does not consider explainability (XAI) to be an add-on requirement.

Model-Agnostic Explainers More recent tools such as SHAP (SHapley Additive exPlanations) and LIME are applied post-hoc to explain individual predictions. To give an early warning, SHAP may indicate what features (e.g., "amount of transaction transferred to Jurisdiction X," "break of weekly pattern") contributed the most to the high risk score and in what direction.

- Interpretable-by-Design Models: Preferably, intrinsically interpretable models such as decision trees are used at individual sub-tasks.
- Adversarial Robustness: Attackers can also strive to poison models or devise ways to misuse transactions in order to cheat. Domain-aware systems run adversarial training and continuous data drift and concept drift (e.g. a change in criminal behaviour) monitoring.
- Model Governance & Audit Trail: A strict governance structure cannot be compromised. This involves version control of the models and features, extensive documentation of model development and validation, ongoing monitoring of its performance based on the key measures (precision, recall, false-positive rate) and a clear audit trail of each alert generated, in connection to the model version and contributing factors.



**Fig 4: Network Risk Propagation Using Graph Neural Networks**

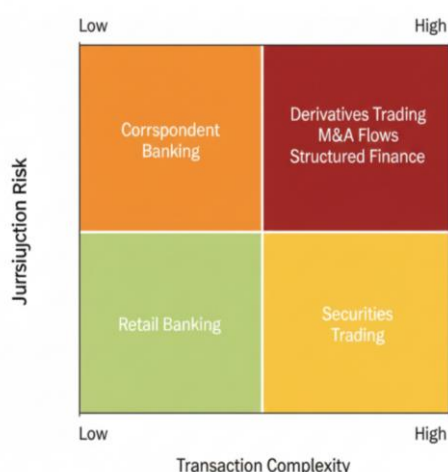## 2.5. Applications in Investment Banking and Cross-Domain Applications

The AML in the domain of investment banking has its own peculiarities.

- Complex Products: The features should code the risk profiles of mergers and acquisitions flows, security trading, and derivative settlements. A massive, non-cumulative collateral posting of a derivatives contract can be normal (Pérez, 2025); the same amount transferred to a personal account in an offshore center, is not. This is a situation that the model should be aware of.
- Client Sophistication: Clients may be individual corporates, financial entities, or ultra-high-net-worth individuals whose legitimate global business

operations are complex. The normal setting of the system should be less strict and strict than in the case of retail banking, but should be stricter in detection of misuse of advanced channels.

- Cross-Border Integration: The system needs to work in harmony with data and regulations among various jurisdictions where the bank is established.
- The domain-conscious ML principles can be applied across the financial services.
- Retail Banking: Pay attention to the process of identifying mule accounts and fraud rings as well as low-value/high-volume laundering.
- Wealth Management: Scanning hidden illicit assets and evasion of sanctions by use of investments.
- Insurance: Fighting fraud and laundering with undisclosed payouts and policy investments.



**Fig 5: Investment Banking AML Contextual Risk Matrix**

## 3. Conclusion

The war against financial crime is an arms race. In a digital, globalized economy, a stagnant, rules based AML system no longer serves the purpose. This paper has posited that the way to go is in a well-considered and carefully planned combination of machine learning and deep compliance domain knowledge- a paradigm we refer to as Domain-Aware Machine Learning.

Financial institutions can significantly revamp their effectiveness by creating intelligent AML systems, the smarts of which are informed by subject-matter expertise across the entire span, such as the formulation of meaningful features and the choice of hybrid model architecture to the utilization of NLP to automate and the relentless dedication to explainability and sound governance. There is a high promise that such systems will bring the efficacy rate of sub-1% SAR detection up enormously, and at the same time, lower the operational expense of false positives. They facilitate a transition to risk-based monitoring rather than reactive monitoring based on thresholds that can adapt to changing strategies of the enemies.

In the case of investment banks and other advanced financial institutions, such a thing is not just a matter of efficiency but it is also a strategy. Intelligent AML systems that are domain-aware will need long-term cooperation between data scientists, compliance experts, and regulators to establish the trust in advanced analytics. The future of financial crime compliance is not to substitute the human judgment with AI but rather to enrich it with intelligent, transparent, and context-driven systems that can allow the investigators to concentrate on what they do best to make critical risk decisions.

## References

1. Hanif, R., Ahmad, H. S., & Ali, A. (2025). Developing an Integrated AML Risk Management Framework for Commercial Banks Based on Customer Risk Profiling and Enhanced Due Diligence. *Advance Journal of Econometrics and Finance*, *3*(3), 206-215. http://ajeaf.com/index.php/Journal/article/view/114

2. Ofoegbu, K. D. O., Osundare, O. S., Ike, C. S., Fakeyede, O. G., & Ige, A. B. (2024). Proactive cyber threat mitigation: Integrating data-driven insights with user-centric security protocols. *Computer Science & IT Research Journal*, *5*(8), 2083-2106. https://www.researchgate.net/profile/Kingsley-Ofoegbu/publication/383606826_Proactive_cyber_threat_mitigation_Integrating_data-driven_insights_with_user-centric_security_protocols/links/66d3512c2390e50b2c21f33b/Proactive-cyber-threat-mitigation-Integrating-data-driven-insights-with-user-centric-security-protocols.pdf

3. Liu, H., Li, Y., & Wang, H. (2025). Genomas: A multi-agent framework for scientific discovery via code-driven gene expression analysis. *arXiv preprint arXiv:2507.21035*. https://arxiv.org/abs/2507.21035

4. Frigiola, A. (2024). *Supervised Contrastive Learning for Classification of Market Stock Series* (Doctoral dissertation, Politecnico di Torino). https://webthesis.biblio.polito.it/31086/

5. Kang, Y., Yang, X., Wang, G., Wang, Y., Wang, Z., & Liu, M. (2025). Can large language models effectively process and execute financial trading instructions?. *Frontiers of Information Technology & Electronic Engineering*, *26*(10), 1832-1846. https://link.springer.com/article/10.1631/FITEE.2500285

6. Yang, H., Zhou, Y., Ji, X., Liu, Z., Tian, Z., Tang, Q., & Shi, Y. (2025). Advancing Graph Neural Networks for Complex Relational Learning: A Multi-Scale Heterogeneity-Aware Framework with Adversarial Robustness and Interpretable Analysis. *Mathematics*, *13*(18), 2956. https://www.mdpi.com/2227-7390/13/18/2956

7. Pérez, J. V., Chávez, M. R., Prieto, M. D., & Martínez, L. R. (2025). Recent progress of anomaly detection in energy applications: a systematic literature review. *Anomaly Detection-Methods, Complexities and Applications: Methods, Complexities and Applications*, 3.

https://books.google.com/books?hl=en&lr=&id=pcKRE
QAAQBAJ&oi=fnd&pg=PA3&dq=Bridging+Data+Sci
ence+and+Compliance:+Intelligent+AML+Systems+by
+Design:+Domain-
Aware+Machine+Learning.&ots=SUN0cJPUsZ&sig=m
E6Ta-_vQ2bHMkBtBzTdto-b0jw

8. Espinoza, M. (2025). *The Effect of Enhancing Medicare Claims With OSINT on the Community Membership of Fraud Rings* (Doctoral dissertation, Marymount University).
https://search.proquest.com/openview/c08256382439ac1
2b574fdfa4b13dd76/1?pq-
origsite=gscholar&cbl=18750&diss=y