**Research Article**

# Neuro-Symbolic Enforcement Engines for Proactive Financial Crime Prevention

## Mallikarjun Reddy Gouni*

University of Illinois Springfield, USA
* **Corresponding Author Email:** gounimallikarjunreddy@gmail.com - **ORCID:** 0000-0002-5007-7590

**Abstract:**

Financial organizations are faced with unprecedented challenges in identifying complex financial crimes that utilize AI, Deepfakes, and Multi-Level Obfuscations. Current compliance solutions are far from adequate by virtue of high levels of false positives and ineptness in spotting new forms of money crimes. Neuro-symbolic enforcement engines can be considered revolutionary solutions that seek to combine neural-based anomaly recognition and symbolic problem-solving capabilities for the proactive prevention of financial crimes. These novel solutions seek to combine transformer-based sequence models for temporal analysis of financial transactions with graph neural networks that represent regulatory policies as symbolic logic structures. These engines enable the system to recognize complex patterns in high-value financial transaction data as well as make rationalized decisions based on formalized compliance rules. Contrastive learning strategies can be used for improved identification of hidden criminal patterns in financial data by adequately addressing the high levels of class imbalance commonly found in Anti-Fraud analytics. Proactive predictive simulation for compliance outcomes on potentially criminal activity before escalation can be used for preemptive action plans. Generative models can be used for simulating new money crime scenarios for adversarial Validation. Real-time processing requirements for enforcement engines and satisfaction conditions for fairness on diverse customer sets can be considered as challenges for implementation.

## 1. Introduction and Problem Context

Financial institutions today are faced with rising challenges of combating smart financial crime, the intelligence of which is rising with increasing expertise in artificial intelligence technology. Modern fake transactions involve deep fake technology for identity forgery, artificial data designs that are undetectable, and laundering structures with multiple layers that cut across geographical and asset jurisdictions. The main challenge with modern compliance technology, primarily a static rules-engine technology, is that it contains inefficiencies of too many false positives and limitations on discovering new types of financial crime that do not fit its histories.

The architectural challenges posed by legacy anti-money laundering solutions come from their deterministic models, which are based on thresholds, geography, and comparisons with customer profiles to predefined templates. These solutions are purely reactive; they mark a transaction solely after determining any suspicious behavior from a predefined suspect pattern, which exists as a curated rule base that a human maintains manually. These approaches are inefficient when dealing with a cunning criminal organization, as they will continue to modify their models according to the existing capability to detect them. A gap exists between capability and criminal sophistication, which results in a high cost of compliance with little efficacy regarding risk mitigation.

Neurosymbolic AI symbolizes an emerging computing paradigm that combines neural net architectural designs and symbolic reasoning platforms to capitalize on the complementary benefits of both. While the combination of neural and symbolic AI allows systems to conduct learning processes dependent on perception from raw input data at the same time as logical reasoning about symbolic forms of knowledge, it fills the major limitations of connecting and symbolic AI methods. Similarly, neural AI methods are highly efficient in neutral application areas such as image classification, modeling, and anomaly detection,

but are not efficient in reasoning and systematically learning general processes in novel settings. On the contrary, symbolic AI methods are highly efficient in logical reasoning and composing knowledge processes, but are not scalable because of significant manual development processes [1].

Using neuro-symbolic architectures for the prevention of financial crimes marks the beginning of a paradigm shift for proactive intervention based on predictive compliance risk analysis instead of the current reactive approach based on the flagging of completed suspicious transactions. These hybrid models typically integrate deep learning networks capable of latent feature extraction from the data stream related to the transaction activity and reasoning engines tasked with assessing the extracted patterns based on formalized frameworks related to the financial regulations. The approach has been applied for the detection of financial crimes related to the flow of cryptocurrencies on graphs, where the graph convolutional neural network has been found effective for the extraction of network-level structural patterns related to money laundering activities [2].

## 2. Theoretical Foundations of Neuro-Symbolic Architectures

Neurosymbolic architectures arise as a result of the appreciation of the fact that neural network methods and symbolic processing systems have complementary strengths, which tackle different problems in intelligent reasoning. Deep neural models uncover hidden regularities in data spaces with higher dimensionality using hierarchical feature discovery, which achieves generalization by example without resorting to rule-based programming. These models accept unstructured data, which may be images, text, or time series data, to name a few, to uncover hidden regularities in the data that can be utilized for predictive tasks. The problem with neural networks is that they act like black boxes with poor explainability, intense training data requirements to attain robustness, and an inability to incorporate hard logical constraints.

Symbolic artificial intelligence encodes knowledge in formal logical structures such as predicate logic, semantic networks, and ontologies that can facilitate deductive inference. Symbolic AI has traceable inferential computations with conclusions that can be traced back to the underlying axioms via documented reasoning chains. Symbolic AI can naturally embed domain knowledge through knowledge engineering. This type of AI is hindered by the lack of autonomous machine learning from examples. Symbolic AI requires an extensive human knowledge base development that is impractical for complicated knowledge domains. It is brittle when faced with incomplete or noisy data examples that differ from the hypothesized regularities.

Abductive learning is a paradigm for bridging neural perception and symbolic reasoning using cycles of refinement where neural models develop explanations for observed facts and symbolic reasoners check hypothesis conformity with knowledge expressed in a symbolic representation. The abductive learning paradigm is for scenarios involving incomplete supervision by learning data and knowledge expressed in symbolic representations. The neural models develop explanations for observed phenomena using learning done on available instances, and logical reasoners check for conformity with logical constraints and inconsistencies to be addressed during hypothesis refinement. In this two-way exchange, there is complementarity in using evidence based on both data and logical knowledge to attain capabilities surpassing those possible using either neural models or logical reasoners [3].

Probabilistic logic programming systems give a mathematical basis for neurosymbolic integration by combining logical reasoning with probabilistic inference on uncertain knowledge. These systems model knowledge using logical predicates, with additional probability distributions modeling uncertainties in facts and rules. Neural networks learn probability values for logical predicates from examples, while logical inference systems reason with uncertainties in rules to compute probability distributions for logical conclusions. Deep ProbLog integrates these technologies by nesting a neural network inside a probabilistic logic program, with outputs of a deep network being used as a probability distribution for logical facts that reason with logical inference engines. The system allows back-propagation of gradients from logical conclusions through logical inference systems to deep network weights, making it efficient for tasks with dual requirements of pattern recognition in raw sensory input and symbolic reasoning over knowledge structures [4].

## 3. Neural Anomaly Detection Framework

The Transformer models yield elementary components that can represent sequential transaction data efficiently through self-attention mechanisms, which enable model awareness of long-term sequential dependencies without recurrent links that cause gradient propagation challenges. The Transformer model is designed to examine input sequences by computing the weight of attention between each sequence position, which

impacts other sequence locations, hence facilitating parallel computation according to different input sequence lengths while considering temporal dependencies. The multi-head attention mechanism enables parallel attention towards different information aspects within distinct representation subspaces, which include amount progression relationships, counterparties' interaction behaviors, and temporal clustering behaviors, among others. The positional embedding techniques allow the model to consider temporal aspects of transaction data by identifying repeated transaction instances that are distinct due to temporal differences [5].

In the case of financial transaction analysis, the transformer encoders deal with a sequence where each transaction is a set of several features, such as amounts, timestamps, identifiers for the counterparties, type information, and information about the accounts. The transformer model can learn the context-related features for each transaction, depending on the context provided by the whole set of behaviors, rather than processing each transaction independently. The attention component is crucial for automatically identifying the set of past transactions that contribute the most to the assessment of current transaction legitimacy.

Contrastive learning techniques train CNNs on embedding spaces where similar data points group together, and dissimilar data points apart, without relying on large quantities of labeled data. In the contrastive learning framework, the CNNs learn to maximize similarities between differently transformed views of the same data point and minimize similarities between differently transformed views belonging to different data points. This technique helps balance the class-imbalance problem in financial crime analysis, where fraudulent transactions only make up very small proportions of overall financial transactions, thereby making supervised learning difficult due to the lack of positive examples [6].

Self-supervised contrastive learning identifies inherent clustering patterns in the transaction data without relying on fraud labels by considering temporal segments of the same customer as positives and segments of different customers as negative samples. The learned embeddings encode fraud behavior consistency patterns common in individual customers, allowing for fraud anomaly detection based on the identification of transactions with irregular behavior compared to predefined customer behavior profiles. Supervised contrastive learning includes fraud samples with real fraud examples by pairing transactions with common fraud properties regardless of underlying specification implementations, promoting the network to learn fraud-related commonalities

despite differing implementations of fraud. This jointly trained approach combines the use of unlabeled transaction data with self-supervised objectives and supervision with fraud labels with supervised objectives, which results in embedding spaces in which novel fraud variants lie close to real fraud examples despite differing on the surface characteristics [6].

## 4. Symbolic Reasoning and Regulatory Alignment

Graph neural networks offer computational models for learning from graph-structured data where the graph nodes represent entities and edges signify relationships between entities. Graph convolutional networks are an extension of graph convolution for irregular graph structures instead of the regular grid structures inherent in convolution. A graph convolutional layer performs a transformation on the nodes of the graph based on the features from the surrounding entities along a particular edge, and a layer in the graph convolutional network updates its features based on the current features and the features obtained from its surroundings for each node [7]. By stacking multiple layers, the features can flow through the graph structure.

Inductive learning over graphs makes it possible to achieve generalization over unseen graph structures through the induction of aggregation functions over sets of node features, thereby avoiding direct encoding of the graph topology in model parameters. This becomes vital in regulatory compliance reasoning, where the knowledge graph keeps expanding over time due to the evolution of new regulations and modifications made to the existing ones. This inductive technique can effectively induce the generation of embeddings of nodes in accordance with their features and neighborhoods, in contrast to direct reference to the identities of nodes, allowing the model to dynamically adapt to an evolving knowledge graph over regulatory concepts without requiring any modifications to be made in the graph topology due to the introduction of new concepts or adjustments made to the relationship among existing concepts [7].

Graph neural network architectures carry out the inference of logical rules through message passing schemes in which activation of connected nodes takes place according to the learned aggregate function. The aggregate function can learn to represent the activation of a node only when its neighboring condition node is satisfied for conjunction, or activation of the node when its neighboring condition node is satisfied for disjunction. The learnability derived from the

differentiable nature of the aggregate function in neural networks steers the learning of the regulatory compliance relationship through training on examples used for assessing such compliance.

The graph neural network model offers theoretical bases for recursive architectures of neural networks dealing with structured data, proving that specifically designed aggregation functions can tackle graph-based learning problems with any level of precision with unlimited computational power. The model analyzes graphs by iteratively updating the state, where each node keeps an updated state vector according to the states of neighboring nodes and their properties until reaching conditions of equilibrium. This model tackles the cyclic dependencies between knowledge requirements, where deregulatory requirements cite other requirements with interdependent relations, which allows complex compliance checks with interdependent phrases to be assessed. This framework is capable of dealing with graphs with different sizes and structures, which allows compliant reasoning between different countries with unique structures under various levels of their regulation structures and during different years where the structure of their regulation frameworks is changed [8].

## 5. Predictive Regulatory Simulation and Synthetic Scenario Generation

Generative adversarial networks provide the foundation for the training of generative networks via adversarial training methods, with the generator network being responsible for the generation of data samples, and the discriminator network being responsible for the discrimination of the generated samples from the true training data. The generator network utilizes random noise vectors as input and applies transformations to the input data to produce output in the data domain, with the discriminator network being responsible for the processing of the true training data and the output from the generator network by assigning probability scores for the input coming from the training distribution versus the generator distribution [9].

This training process between adversaries causes the generator to keep generating more realistic data point replicas, as the generator with outputs detectable by the discriminator receives a strong learning signal, while the one with believable outputs receives a weaker learning signal that implies successful tricking of the discriminator. The generator ultimately attains equilibrium, learning to produce data points that cannot be distinguished from the training data by the learned standards of the discriminator, thus learning to sample from the

training data distribution. The conditional generative adversarial networks are an extension of this process that condition the generator and its corresponding discriminator on class labels, among others, thus facilitating the generation process according to the required attributes.

In respect of the task of generating financial crime scenarios, conditional adversarial networks can generate a series of transactions with defined characteristics of a certain pattern of fraud by conditioning on labels of fraud types and specifications of parameters. In this manner, a generator learns to map defined noise distributions, along with other conditions, to realistic transactions with a specified pattern of fraud, emphasizing the incorporation of features of actual transactions, including defined amounts, counterparty relationships, and other features that fit certain specified labels of a certain pattern of fraud. The generated scenarios serve for testing compliance engine effectiveness in a broad range of different manifestations of fraud without actually conducting tests on actual fraudulent transactions, allowing for an analysis of compliance engine weak points [9].

The large language models are trained on extremely large corpora, achieving the capability for "few-shot learning," in which the models can accomplish the task with very few training examples in the task itself while using the general knowledge gained while pre-training the model. The model has shown the capability to accomplish the task for "emergent reasoning, long-form coherent text generation, and natural language directives defining the task steps" using the natural language prompts via the "few-shot learning" paradigm, which specifies the task using example input-output pairs, requiring no gradient-based fine-tuning on the task data [10]. In the context of synthetic fraud scenario generation, language models can generate narratives about possible fraud scenarios based on conditioning their generation on certain fraud characteristics, vulnerabilities, and limitations. These narratives, therefore, embody the creativity in fraud scenarios based on known fraudulent patterns and showcase new attack vectors that combine known methods in a completely new way. These narratives are rendered into parameterized specifications through extraction processes that convert the narratives into parameterized representations suitable for rule-based generation engines that can generate transaction sequences [10].

## 6. Implementation Considerations and Operational Integration

The variational autoencoder offers probability models involving the training of encoder networks,

which map inputs to probability distributions of the latent codes, and decoder networks, which map samples from these probability distributions to inputs. The training process is balanced between the reconstruction of inputs and the regularization of the latent distributions approaching specified prior distributions, like the standard Gaussian distribution, that allow easy sampling and interpolation within the latent space, thus permitting the generation of new data by sampling the latent codes from these prior distributions [11], which is useful in anomaly detection via error analysis of reconstructed inputs that show deviational data points.

In transaction risk score modeling, variational autoencoders learn to define a low-dimensional latent representation that encodes key behavioral traits that characterize transactions conducted by each specific customer. Such latent representations are derived using an encoder that maps a set of transactions to define latent distributions that correspond to behavioral representations of different customers, with a decoder that maps latent examples to transactions. Transactions that report high reconstruction error values are those that do not follow regular behavioral representations embedded in latent space, acting as anomaly indicators complementing those identified using transformer sequence modeling techniques.

The variational autoencoding framework easily accommodates varying attribute types found in transactions using suitable designs for the encoder and decoder functions operating with a mix of data types, which include numerical, identifier, and timestamp data types, respectively. The embedding space found with these models offers a common view across varying attribute types, allowing for a global assessment of user behavior with insights culled from a collection of data types of varying characteristics. Variational architectures designed with a hierarchy for modeling user behavior over varying time scales exploit level-wise embeddings

where top-level embeddings encode higher, longer-term user behavior, and bottom-level embeddings encode lower, contemporaneous behavior changes necessary for risk assessment [11].

Algorithmic fairness relates to biased outcomes related to protected demographic groups regarding disparate impact or discrimination. Algorithm design considering fairness involves constraints that set demands related to statistical parity, like demographic parity, where acceptance rates are equal across groups, or equalized odds, where the true positive and false positive rates are equal across groups. In contrast, there are considerations related to accuracy, where optimal models often demonstrate disparate outcomes across groups when base rates differ across groups according to demographic attributes. Algorithm design considering social aspects involves several objectives, such as accuracy, fairness in relation to groups, or algorithmic interpretability related to accountability in algorithmic decisions [12].

Financial fraud detection engines face challenges balancing achieving the greatest possible effectiveness in fraud detection with avoiding unfairly biased treatment of client groups with different demographic profiles. Fraud risk-assessing models learning from experience with biased client demographic profiles may pick up on demographic correlates indicating high fraud risk associated with disproportionate fraud enforcement in the past, effectively reinforcing biased outcomes with automated fraud-detection engines. Fairness during model development involves adding constraints that demote disproportionate outcomes for protected classes during model optimization, leading to risk scores meeting given fairness requirements without sacrificing fraud-detection effectiveness. Symbolic reasoning parts with interpretability capabilities facilitate scrutiny of fraud-detection engine logic for possible demographic discrimination in accordance with anti-discrimination legislation protecting users in the financial industry [12].

*Table 1: Neurosymbolic Integration Approaches and Their Characteristics [3, 4]*

| Integration Paradigm | Learning Mechanism | Reasoning Capability | Primary Application Domain | Knowledge Representation |
|---|---|---|---|---|
| Abductive Learning | Hypothesis generation from examples | Constraint verification and consistency checking | Incomplete supervision scenarios | Symbolic predicates with neural perception |
| Probabilistic Logic Programming | Probability assignment for predicates | Uncertainty propagation through rules | Pattern recognition with logical inference | Probabilistic facts grounding logic programs |
| Pure Neural Networks | Hierarchical feature discovery | Limited compositional generalization | Image classification and sequence modeling | Distributed representations |
| Pure Symbolic Systems | Manual knowledge engineering | Deductive inference over axioms | Formal reasoning domains | Predicate logic and ontologies |

***Table 2:*** *Neural Architecture Components for Transaction Anomaly Detection [5, 6]*

| Architecture Component | Computational Mechanism | Temporal Dependency Handling | Feature Learning Strategy | Class Imbalance Mitigation |
|---|---|---|---|---|
| Transformer Encoders | Self-attention weights across sequence positions | Positional encoding for temporal context | Contextualized transaction representations | Multi-head attention for diverse patterns |
| Self-Supervised Contrastive Learning | Maximize similarity within customer segments | Temporal segments as positive pairs | Behavioral consistency embeddings | Unlabeled data utilization |
| Supervised Contrastive Learning | Maximize similarity across fraud instances | Cross-instance fraud pattern clustering | Fraud-indicative feature generalization | Limited labeled example leverage |
| Recurrent Networks | Sequential hidden state propagation | Built-in temporal modeling | Gradient-based sequence learning | Limited gradient propagation |

***Table 3:*** *Graph Neural Network Capabilities for Regulatory Reasoning [7, 8]*

| GNN Capability | Graph Processing Method | Adaptability Mechanism | Logical Operation Approximation | Structural Flexibility |
|---|---|---|---|---|
| Inductive Learning | Aggregation over node feature sets | Generalization to unseen graph structures | Feature-based node embedding generation | Dynamic regulatory concept addition |
| Message Passing | Information exchange between connected nodes | Learned aggregation functions | Conjunction and disjunction operations | Multi-layer information propagation |
| Recursive State Updates | Iterative node state refinement | Equilibrium-based convergence | Cyclic dependency resolution | Variable graph topology handling |
| Graph Convolution | Neighbor feature aggregation | Layer-wise representation transformation | Structured pattern capture | Extended neighborhood analysis |

***Table 4:*** *Generative Model Applications for Fraud Scenario Synthesis [9, 10]*

| Generative Model Type | Training Mechanism | Synthesis Control Method | Output Characteristics | Validation Purpose |
|---|---|---|---|---|
| Generative Adversarial Networks | Adversarial discriminator-generator optimization | Conditional fraud typology labels | Realistic transaction sequences with fraud patterns | Detection blind spot identification |
| Conditional GANs | Class-conditioned generation process | Parameter specification for fraud attributes | Controlled fraud characteristic manifestation | Adversarial robustness testing |
| Large Language Models | Few-shot learning from prompts | Natural language fraud specifications | Creative fraud scenario narratives | Novel attack vector exploration |
| Rule-Based Synthesis | Parametric specification translation | Structured extraction from narratives | Transaction sequences with statistical realism | Compliance engine effectiveness evaluation |

## 7. Conclusions

Neuro-symbolic enforcement engines signify a paradigm shift in financial crime prevention, combining pattern recognition employing neural networks and logical reasoning using symbols to effectively compensate for the inadequacies of traditional compliance systems. The neuro-symbolic approach enables parallel learning from data streams and systematic reasoning based on formalized regulatory systems to provide interpretable compliance evaluations based on legitimate logical inference paths. Transformer models provide sequence-level dependencies for complex transactions over extended timeframes, and contrastive learning enables generalized embeddings that are effective across varied forms of fraud occurrences. Graph neural networks are utilized to implement regulatory graphs that facilitate automated testing of identified violations based on associated compliance laws and regulations within specific geographical jurisdictions. Generative adversarial networks and large language models are applied to develop

imitation fraud holistically, allowing systematic vulnerability analysis for regulatory improvement without actual occurrences of fraud. Variational autoencoders introduce probabilistic measures for risk assessment, employing uncertainty quantification measures for behavioral analysis. The complementary technologies combined are effective in setting up enforcement systems beyond specific reactive transaction flagging systems that deliver enforcement capacity for closely proactive measures based on predictive compliance risk analysis. Fairness-oriented design strategies are applied to ensure balanced treatment of various demographic groups while preserving effective fraud detection performance. Such milestones position neuro-symbolic technologies to be at the root of next-generation financial fraud protection systems sensitive to drastically evolving fraudulent practices.

## Author Statements:

- **Ethical approval:** The conducted research is not related to either human or animal use.
- **Conflict of interest:** The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper
- **Acknowledgement:** The authors declare that they have nobody or no-company to acknowledge.
- **Author contributions:** The authors declare that they have equal right on this paper.
- **Funding information:** The authors declare that there is no funding to be acknowledged.
- **Data availability statement:** The data that support the findings of this study are available on request from the corresponding author. The data are not publicly available due to privacy or ethical restrictions.

## References

[1] Artur d'Avila Garcez and Luis C. Lamb, "Neurosymbolic AI: The 3rd Wave," arXiv:2012.05876, 2020. [Online]. Available: https://arxiv.org/abs/2012.05876

[2] Mark Weber et al., "Anti-Money Laundering in Bitcoin: Experimenting with Graph Convolutional Networks for Financial Forensics," arXiv:1908.02591, 2019. [Online]. Available: https://arxiv.org/abs/1908.02591

[3] Wang-Zhou Dai et al., "Bridging Machine Learning and Logical Reasoning by Abductive Learning," 33rd Conference on Neural Information Processing Systems, 2019. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/20 19/file/9c19a2aa1d84e04b0bd4bc888792bd1e-Paper.pdf

[4] Robin Manhaeve et al., "DeepProbLog: Neural Probabilistic Logic Programming," arXiv:1805.10872, 2018. [Online]. Available: https://arxiv.org/abs/1805.10872

[5] Ashish Vaswani et al., "Attention Is All You Need," arXiv:1706.03762, 2023. [Online]. Available: https://arxiv.org/abs/1706.03762

[6] Ting Chen et al., "A Simple Framework for Contrastive Learning of Visual Representations," arXiv:2002.05709, 2020. [Online]. Available: https://arxiv.org/abs/2002.05709

[7] William L. Hamilton, Rex Ying, and Jure Leskovec, "Inductive Representation Learning on Large Graphs," arXiv:1706.02216, 2018. [Online]. Available: https://arxiv.org/abs/1706.02216

[8] Franco Scarselli et al., "The Graph Neural Network Model," IEEE Transactions on Neural Networks, Volume 20, Issue 1, 2009. [Online]. Available: https://ieeexplore.ieee.org/document/4700287

[9] Ian J. Goodfellow et al., "Generative Adversarial Networks," arXiv:1406.2661, 2014. [Online]. Available: https://arxiv.org/abs/1406.2661

[10] Tom B. Brown et al., "Language Models are Few-Shot Learners," arXiv:2005.14165, 2020. [Online]. Available: https://arxiv.org/abs/2005.14165

[11] Diederik P Kingma and Max Welling, "Auto-Encoding Variational Bayes," arXiv:1312.6114, 2022. [Online]. Available: https://arxiv.org/abs/1312.6114

[12] Michael Kearns and Aaron Roth, "The Ethical Algorithm: The Science of Socially Aware Algorithm Design," Oxford University Press, 2019. [Online]. Available: https://dl.acm.org/doi/book/10.5555/3379082