

Self-Directed Intelligence in Finance: Agentic AI Models for Cross-Border Risk and Compliance Automation

Pallav Kumar Kaulwar

IT Director, pallavvkumar@gmail.com, ORCID ID: 0009-0002-1142-0329

ARTICLE INFO

Received: 02 Nov 2024

Revised: 16 Dec 2024

Accepted: 26 Dec 2024

ABSTRACT

Intelligent Systems and Advanced Scientific Computing In finance—in fact, in other industries as well—huge amounts of data are collected and processed as a part of operations. With self-directed intelligence capabilities, these operations in finance can lead to a very low cost structure, especially in the case of risk and compliance checks. A cost structure that is much lower than the revenue generated by providing services will lead to very profitable financial institutions. Self-Directed Intelligence—Self-Directed Intelligence (SDI) is a type of intelligence that is independent of human intervention. In finance, SDI is used to detect anomalies or patterns in transactions. Rather than waiting for alerts from other systems, a self-directed intelligence system continuously ingests large amounts of transaction data and generates workflows, reporting or other operational tasks autonomously. SDI is a major pillar of Artificial Intelligence. Artificial Intelligence (AI) is the popular term today for intelligent systems with self-directed intelligence for specific use cases defined for self-directed systems. Self-Directed Systems (SDS) is another umbrella term covering both SDI and AI that includes systems with a lower level of autonomy supported by human experts at all times. Agentic AI focuses on the systems with the highest level of autonomy—agentic and self-directed—that operate independent of human intervention. For the applications described in this paper—detecting anomalies or patterns in transactions and generating workflows, reporting, or other operational tasks autonomously—AI with a low operating threshold is sufficient. Agentic AI is a layered stack that can operate at different levels of autonomy.

Keywords: Self-directed intelligence, Agentic AI systems, Autonomous financial agents, Cross-border financial compliance, AI-driven risk management, Regulatory technology (RegTech), Automated AML monitoring, AI for KYC automation, Cross-jurisdictional risk assessment, Intelligent compliance orchestration, Financial governance automation, Multi-regulatory AI frameworks, Adaptive compliance systems, Real-time regulatory monitoring, AI-powered fraud detection, Autonomous decision-making in finance, Compliance-by-design AI, Global financial risk analytics, Explainable AI for compliance, Policy-aware AI agents.

1. Introduction

Globalization, technological change, and regulatory complexity amplify the risk of finance becoming a localized service industry reliant on transborder data flows. Yet, self-directed intelligence affords the capability to automate service provision across organizational and geographical boundaries. Cross-border risks arise from asynchronous economic cycles, monetary and credit policies, counterparty exposures, and deviations in market, operational, or location risk considerations. Such risks have long been recognized; the 2008–2010 financial crisis and, more strongly, the COVID-19 pandemic have brought them back into sharp focus. Deleterious cross-border effects have been mitigated by statistics-driven macro-prudential measures, but the risk of liquidity bottlenecks in times of stress remain. Alongside the interdependencies of financial markets and institutions, supervision and regulatory philosophies have also diverged—reflected in different approaches to financial and crypto-asset market infrastructure.

The ongoing far-reaching geopolitical rift—as channeled through the decoupling of China and the United States—poses profound challenges for major financial institutions worldwide, creating a new and complex range of risks. In parallel, there has been a surge in the number of criminal and terrorist cases involving illicit financing by financial institutions, sparking the need for agencies and authorities to undertake continuous and more concerned anti-money laundering (AML) monitoring to satisfy expanding legal, fiscal, and regulatory requirements. AML-related compliance has been recognized as one of the fastest-growing and most resource-intensive legal compliance areas; institutions are looking to technology to drive greater efficiencies and effectiveness within their operations.

The widening geopolitical divide—most notably reflected in the strategic decoupling between China and the United States—has introduced a heightened level of uncertainty and systemic risk for global financial institutions, complicating cross-border operations, regulatory alignment, and risk management frameworks. At the same time, the growing prevalence of criminal and terrorist financing cases linked to financial institutions has intensified scrutiny from regulators and enforcement agencies worldwide. This has driven the demand for continuous, more sophisticated anti-money laundering (AML) monitoring to meet increasingly stringent legal, fiscal, and regulatory obligations. As AML compliance has emerged as one of the fastest-growing and most resource-intensive areas of regulatory compliance, financial institutions are increasingly turning to advanced technologies—such as automation, data analytics, and artificial intelligence—to enhance detection capabilities, improve operational efficiency, and strengthen the overall effectiveness of their compliance programs.



Fig 1: Finance: Agentic AI

1.1. Background and Significance

Self-directed intelligence is the capability of an agent to identify and pursue its own goals through a closed decision-making and action-taking loop. Agentic AI systems realize self-directed intelligence on behalf of their users by combining elements from building blocks such as decisioning cycles, data processing stacks, learning loops, decisioning cycles, task loops, and human-in-the-loop interfaces. The concept of autonomy encompasses the level of independence with respect to the user and the boundary of governance around agentic capabilities. In Finance, agentic AI models may assume responsibility for specific risk and compliance requirements across organizational and geographical boundaries.

Regulatory authorities in Finance have become increasingly concerned about the systemic risks to financial stability posed by cross-border activities of banks and non-bank firms. Such risks arise from imperfectly correlated financial shocks across jurisdictions that have distinct economic cycles and regulatory frameworks. The cross-border risk landscape comprises market risk, counterparty risk, operational risk, and compliance risk. Cross-border activities expose firms to compliance challenges because of regulatory heterogeneity, country-specific demands, and the absence of equivalence arrangements. These challenges affect firms' risk and compliance management systems and require the implementation of persistently operating, real-time surveillance, monitoring, and reporting capabilities. To fully address increasingly complex cross-border needs, financial institutions are looking to conduct specific risk and compliance functions through resource-light and cost-efficient systems.

Equation 1: Real-time monitoring & anomaly detection (thresholds)

Let x_1, x_2, \dots, x_N be a monitored metric (e.g., transaction volume index, complaint count, AML alerts/hour).

Step 1 – mean

$$\mu = \frac{1}{N} \sum_{t=1}^N x_t$$

Step 2 – sample variance

$$s^2 = \frac{1}{N-1} \sum_{t=1}^N (x_t - \mu)^2$$

Step 3 – standard deviation

$$\sigma = \sqrt{s^2}$$

1.2 Convert the profile into incident-triggering thresholds

Pick a sensitivity multiplier k (common choices: 2, 3).

Upper threshold

$$U = \mu + k\sigma$$

Lower threshold

$$L = \mu - k\sigma$$

Incident rule

$$\text{Flag}(t) = \begin{cases} 1 & \text{if } x_t > U \text{ or } x_t < L \\ 0 & \text{otherwise} \end{cases}$$

That “Flag” is exactly what feeds a case/incident workflow (as described).

1.3 Useful derived quantities

Z-score (how many sigmas away)

$$z_t = \frac{x_t - \mu}{\sigma}$$

Then $|z_t| > k$ is equivalent to the threshold test.

EWMA (smother real-time baseline)

$$s_t = \lambda x_t + (1 - \lambda)s_{t-1}, \quad 0 < \lambda \leq 1$$

This replaces μ with a moving baseline when the environment drifts.

1.2. Research designs

The study employs a design-oriented approach, supported by secondary sources. Academic literature and industry publications inform a conceptual analysis of self-directed intelligence in finance, agentic AI system components, cross-border risk typologies, compliance automation capabilities, and risk management frameworks. Guidance from ChatGPT, a prototype Agentic AI model, and banking domain expertise provide practical context. Engineering principles and trading guidelines facilitate the development of a trading model, while permutations of self-directed intelligence enable experimentation with a self-learning advisor. Data from crypto trading and multiple assets shape the advisor model.

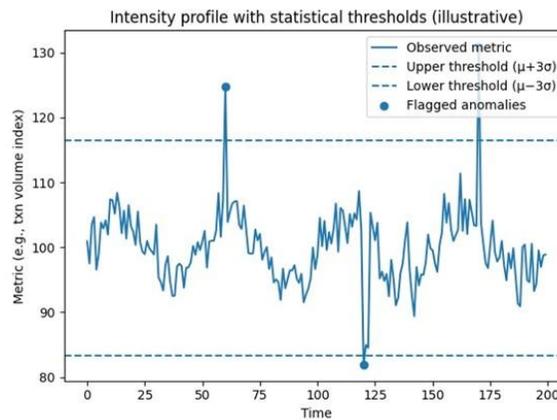
Cross-border risk comprises jurisdictional factors that affect financial institutions, their clients, and transactions. The analysis highlights the heterogeneity of regulatory regimes, the associated demands for compliance, the market, counterpart, and operational risk elements that arise in cross-border activities, and the data sovereignty and privacy concerns that impact data flow. Cross-border regulatory diversity creates challenges for the operations of financial institutions and their clients involved in multi-country cross-border business or trading. Persistent trade imbalances for certain countries attract illegal cross-border capital flows, money laundering activities, and fraud. Counterparty risk is further elevated by higher interconnectedness and systemic importance among major institutions.

2. Conceptual Foundations of Self-Directed Intelligence in Finance

Self-Directed Intelligence in Finance encompasses self-directed applications that fuse supervisory machine intelligence for learning and judgment with humans for strategy and governance. The financial domain supports all financing needs and processes, including risk and compliance. In principle, intelligent automation reduces the need for organizational compliance and risk-control technology capable of autonomous detection and response. However, despite advances using AI for data processing and the rapid evolution of risk-management capabilities, no intelligent applications are able to continuously monitor the regulatory environment, assess compliance with diverse rules, and automatically detect and respond to breaches. In particular, the combined advances of self-learning capabilities and machine-based regulatory

environment monitoring have not yet been harnessed. Therefore, compliance support and guidance from supervisory agents for risk-control organizations are still necessary.

Agentic AI denotes a self-learning architecture that features supervisory intelligence capable of autonomous learning and decision-making within defined boundaries. The architecture consists of a learning loop responsible for processing experiences and a decision-making cycle based on state-reward pairs. The distinctive traits of agentic systems include the operational implementation of autonomy levels beyond full automation, intent specification expressed in natural language, transparency at multiple layers accommodating different consumer needs, and built-in capability for broad-spectrum machine ethics with continuously updated safety profiles. Although data, modeling, and applications have been articulated and analyzed in detail, agentic systems offer complementary potential for finance, where compliance and risk-control organizations remain essential.



2.1. Definitions and Scope

Self-directed intelligence refers to autonomous or semi-autonomous systems able to perform complex and safety-critical tasks without human intervention, adapt to novel situations via learning, and generate their own objectives within a defined governance framework. The distinguishing trait of self-directed intelligence is an internal goal generator: a model processing external stimuli and internal predictions to infer high-level goals. Creating self-directed machines poses complex scientific and technical challenges, but solutions are emerging, primarily within the AI community.

Linguistic conventions have recently changed: self-directed systems are usually referred to as agentic or autonomous AI. Self-directed systems with human-like capabilities are sometimes denoted artificial general intelligence. Agentic AI can describe any technical system or method able to modify its operating characteristics. Autonomous denotes a well-defined operational interface in which accountability remains with one overarching human agency, such as the pilot-in-command of a commercial aircraft.

2.2. Agentic AI: Architecture and Operational Principles An agentic AI system comprises a semantic knowledge base, a long-term memory, a real-time data monitoring facility, an external execution engine, an intelligent agent, and an external environment. At its core lies a semantic knowledge base that provides the meanings of terms in its inputs and outputs by establishing relationships among concepts in the monitored milieu. Contractual relationships and other pacts that specify the contents of expected

communications may also be stored. There exist long-term memory banks that store the lessons learned in the past (including the patterns of interactions with reliable counterparts). In addition, real-time monitoring capabilities analyse streams of data that are currently flowing, comparing them with known patterns to identify possible deviations that may indicate an anomaly.

Three loops form the operational engine of an agentic AI system: a learning loop that improves future performance, a decision-making cycle that makes choices according to the specifications of the agent's intent, and a monitoring loop that keeps an eye on the environment through the data feeds connected to it, detecting anomalies when they occur. The learning loop improves future decisions by feeding back into the system the assessments of the appropriateness of past decisions. Information fed back from the execution engine regarding the results of past actions is used to improve the performance of future decisions down to the decision-making level. The detected anomalies trigger alerts of varying degrees of urgency levels for timely human intervention when required but that are too urgent to wait for the next decision-making cycle. At the broadest level of the architecture, the agent is deemed as truly 'agentic' and fully autonomous when unused because the governance mechanisms do not require immediate human involvement. Such a self-directed capability encompasses all the operative aspects, from continuous data monitoring and situation assessment to detection of anomalies or events that require human monitoring.

2.3. Autonomy, Intent, and Governance in Financial Contexts

Autonomy is a core property of self-directed intelligence. In simple terms, a self-directed system executes a course of action without needing instruction from a human. Such autonomy, however, does not imply intent; systems without intent operate solely on internal programming, following predetermined logic without the ability to self-modify. For example, traditional software applications are not self-directed, for they execute according to explicit instructions coded by humans and thus lack a development feedback loop.

The concept of intent is useful for clarifying a system's operational boundaries – à la – the Agent Design World Model (Davis 2020). The Agent Design World Model proposes deliberate specification of a system's intent for three reasons. First, risks that fall outside the stated intent cannot be effectively controlled. Second, intent alignment during the development cycle becomes impractical if intent is not defined and understood. Third, end-users require assurances that the intended purpose of an agentic AI system is achievable. A system's intent shapes the guardrails of autonomous behavior. With rigorous intent specification and effective intent alignment, any level of action autonomy may be acceptable. The specification of boundaries of action autonomy, however, is distinct from the deliberate delimitation of intent, which guides self-directed systems toward beneficial outcomes.

Equation 2: AML/Fraud “risk scoring” via supervised learning

Let each transaction i have features:

$$\mathbf{x}_i = [\text{amount, frequency, purpose, counterparty-risk, ...}]$$

and label $y_i \in \{0,1\}$ (e.g., suspicious vs non-suspicious).

Step 1 – linear score

$$s_i = \beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i$$

Step 2 – squash to probability

$$p_i = \sigma(s_i) = \frac{1}{1 + e^{-s_i}}$$

Step 3 – operational decision thresholds

$$\text{Decision}(i) = \begin{cases} \text{Auto-approve} & p_i < \tau_1 \\ \text{Escalate/case} & \tau_1 \leq p_i < \tau_2 \\ \text{Auto-reject / block} & p_i \geq \tau_2 \end{cases}$$

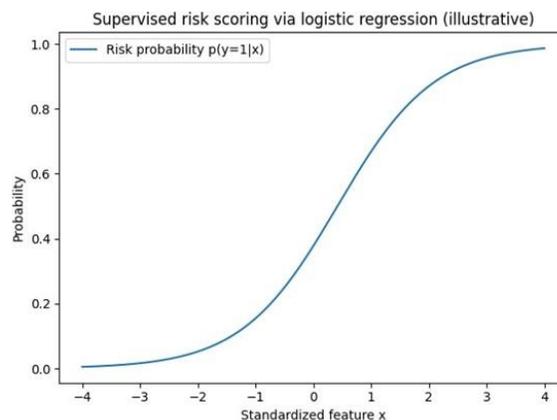
This matches the paper’s “approval, rejection, or escalation”.

3. Cross-Border Risk Landscape

Cross-border financial operations invariably expose institutions to market risk, counterparty risk, and operational risk. Although these risk types are traditionally examined in isolation, they must be considered jointly in global finance because they are intertwined and mutually dependent.

Regulatory frameworks are rarely consistent across jurisdictions. Even where differences are relatively mild, institutions must reconcile divergent rules. In addition to the accumulated cost of compliance, large-scale differences in regulatory approach add another layer of complexity: the danger of incurring sanctions for breaking the rules of the jurisdiction of the regulator of last resort. Such risk is especially pertinent when taking into account conditions imposed on banks under emergency liquidity assistance. It is thus unsurprising that the timely design and implementation of compliance solutions capable of handling these complexities is of paramount importance.

Moreover, the prevailing approach to compliance is fast becoming inadequate. Regulatory authorities no longer consider compliance to be a mere ‘check-the-box’ exercise. Institutions are now expected to have the capability to monitor compliance breaches in real time and immediately report violations to regulators, providing the evidence demonstrating that all measures required to reduce potential damages are being enforced at all times. These levels of assurance cannot be achieved with the traditional, largely manual approach to compliance; much of which still relies on a simple control and reporting cycle with no real-time monitoring or early reporting capability. Such anachronistic procedures need to be replaced as a matter of urgency.



3.1. Regulatory Heterogeneity and Compliance Demands Variation across the legal frameworks of different jurisdictions manifests itself most clearly in anti-money laundering (AML) and anti-terrorist financing (ATF) rules, capital requirements, solvency standards, treatment of credit ratings agencies, business conduct, and conflict of interest policies. The relative prioritization of these issues by different regulators also affects business firms, their cross-border risk exposures, and their economic performance. Compliance with these diverse regulations often requires banks to maintain separate records, procedures, and systems for different countries and markets, leading to inefficiencies and loss of flexibility. Within this context, the availability of a foreign partner regulated in an equivalent manner in the cross-border stakeholders' country clearly constitutes a lower compliance burden, at least for market and counterparty risk. The lack of international harmonization is therefore a constraint for global financial markets. In the immediate future, some sort of equivalence regime is likely to apply to the application of mandatory financial market rules, while cross-border transactional exposure to market and counterparty risk is being dealt with on a more ad-hoc basis. At the same time, data location and privacy will remain very important issues, with strong requirements for data to be stored and processed locally, restrictions on cross-border data transfers, and privacy rules covering data usage. The need to obtain permission from local authorities for such transfers adds to the overhead of cross-border transactions, especially on a one-off basis. The collection of customer data thus introduces operational risk in transactions with certain jurisdictions. The same considerations apply to any communications with an established link or network that traverses a jurisdiction with restrictive data location or privacy rules. In particular, the introduction of a new data location rule can make existing transactions non-viable and appear as an operational risk.

3.2. Market Risk, Counterparty Risk, and Operational Risk in Global Finance Market risk, counterparty risk, and operational risk arise from the very nature of the financial business model. Since transactions involve both buyer and seller, the potential for loss exists in the interest rate, foreign exchange, and commodities markets as a result of market movements. The exposure and consequent risk from changes in FX, interest rate, and commodity prices, together with the risk arising from trading, investment, and commercial activities, usually fall within market risk. A financial institution incurs counterparty risk when an individual trade transaction exposes it to a loss from the default of the counterparty. Operational risk captures the losses that stem from human error, sub-optimal execution or behavior, systems failure, external events, and the general inadequacy of controls.

Losses from these three risk categories interact in several ways. For market risk, the greatest direct impact is on the likely severity of loss over a defined time horizon based on stress testing, limits, reserves, and capital. An absence of active management increases vulnerability. During substantial market turbulence, the probabilities of default (PDs) on credit exposures rise sharply, often to very high levels. Accurate predictions of migration originate primarily from an informed view of credit market conditions. Well-structured and documented loan covenants can help identify potential problems at an early stage. For operational risk, large trading operations can stress many of the elements of an institution's operational risk management. History suggests that the possibility of severe loss from operational risk should always be considered. Stress testing and control function capabilities are important mitigation tools.

3.3. Data Sovereignty and Privacy Considerations The key policies that may affect the deployment of AI and analytics solutions based on the processing of end-user data are: (i) cross-border data flows, (ii) data localization rules, and (iii) personal privacy. Cross-border data flows refer to cross-border physical data flows. In abstract terms, these policies regulate aspects related to the transfer of data

associated with the conduct of economic activities across jurisdictions. Companies need to comply with these requirements when processing data in a foreign jurisdiction or transfer data from or to a foreign jurisdiction. Action plans for the free flow of data, such as those defined by the European Union and Japan, can facilitate an equivalent approach to cross-border data flows and contribute to cross-border data transfer across different jurisdictions with digital economy agreements.

Data localization rules mandate data generated or collected in a given jurisdiction to be stored or processed on servers physically located in that jurisdiction. Data localization goes beyond the protection of personal data and may originate from the protection of interests related to national sovereignty and security. Many countries enact data localization rules covering sensitive personal, financial, and health data in reaction to the unauthorized transfer and use of sensitive national data. Data localization rules can slow down the globalization process of big data solutions based on abundant and diverse data volumes or interviewer data flows. Countries maintaining a similar legal-oriented framework for privacy protection, such as the General Data Protection Regulation, have fewer or no localization rules or follow a risk-based approach to data localization. Data localization rules are only feasible for a limited period of time and should gradually evolve to minimize administrative and business burden.

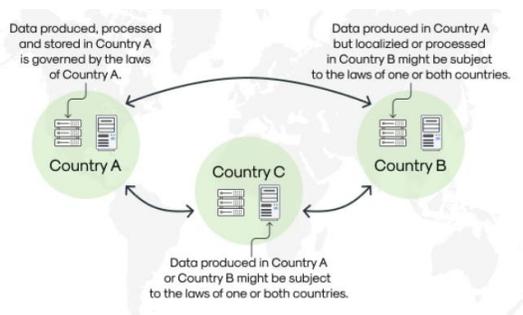


Fig 2: Data Sovereignty

4. Agentic AI for Compliance Automation

Automated compliance systems encompass various functions that can be feasibly executed with existing technology. The automation of two-dimensional tasks with straightforward decision rules requires only basic evidence and testing capabilities that can already be delivered by a basic rule-based engine. In a compliance space shielded from extreme disruptions, automated systems simply need to operate as formally required: monitoring, generating reports and alerts, and making timely decisions.

Market regulators oversee global financial institutions from different jurisdictions. Authorities often require institutions to ensure that foreign branches and subsidiaries comply with local laws. Institutions without a presence in certain jurisdictions must follow rules imposed on third-country institutions. Therefore, business partners face cross-border market risk, counterparty risk, operational risk, and data protection and privacy risk in a surveilled but unshielded space that simultaneously spans multiple geopolitical arenas.

Compliance automation capabilities encompass continuous monitoring for anomalous conduct, automatic generation of regulatory reports, detection of money laundering and fraud, and screening for sanctions compliance. Although potentially highly sensitive and thus traditionally reserved for manual fulfilment,

compliance activities can be efficiently delivered by agentic AI. Fully developed functional capabilities enable an institution to arrange any compliance task within one of these redundant pathways. Each pathway’s fulfilment then guarantees timely closure of the original task.

Equation 3: Agentic “learning loop” and “decision-making cycle” (state–reward)

A standard formalization is a Markov Decision Process (MDP):

- state $s \in \mathcal{S}$ (current compliance/risk context),
- action $a \in \mathcal{A}$ (escalate, request docs, block, report, etc.),
- reward $r = R(s, a)$ (utility: risk reduced, SLA met, false positives penalized),
- transition $P(s'|s, a)$,
- discount $\gamma \in [0,1)$.

3.1 Bellman optimality (decisioning objective)

Step 1 – define value under optimal behavior

$$V^*(s) = \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right]$$

Step 2 – one-step lookahead decomposition

From state s , take action a , get immediate reward $R(s, a)$, then move to s' :

$$V^*(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s'|s, a) V^*(s') \right)$$

That’s the canonical “closed loop” decision+feedback system the paper describes.

3.2 Q-learning update (learning loop)

Define action-value $Q(s, a)$.

Step 1 – TD target

$$\text{target} = r + \gamma \max_{a'} Q(s', a')$$

Step 2 – TD error

$$\delta = \text{target} - Q(s, a)$$

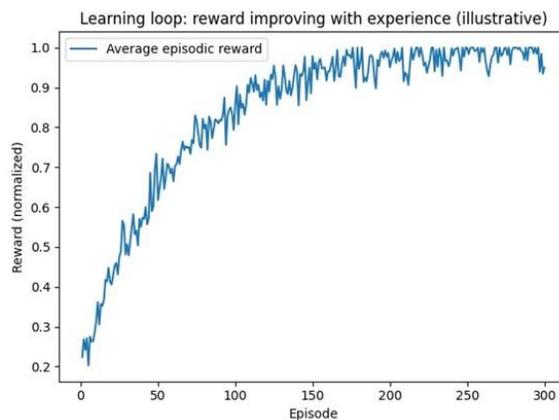
Step 3 – update

$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta$$

(α is the learning rate.)

4.1. Real-Time Monitoring and Anomaly Detection Continuous surveillance of business environments is critical for understanding both normal operation and adverse situations. Most organizations establish intensity profiles around important metrics, like volumes of transactions, customer complaints, counterfeit goods found on the market, or money laundering cases identified. These profiles can be created using traditional statistical techniques or machine learning. The result allows stakeholders to establish both upper and lower thresholds. If one of these thresholds is crossed, a specific incident is triggered. Such incidents can be fed into a case management system, so that a dedicated team investigates the situation. If an anomaly can be confirmed, further investigation can take place, possibly leading to a cascade of related inquiries. Anticipating and planning responses to several categories of these events allows organizations to minimize the potential damage associated with such ruptures.

Every operation running in real-time requires safeguards and must be continuously tested. These tests can also include examining why a particular result was flagged but not investigated in the past. The reasons detected will lead to the creation of new case modes or examination of pre-existing modes to check why those alerts were not acting upon.



4.2. Automated Regulatory Reporting The automated generation of regulatory reports is a straightforward application of agentic AI. Internal clients specify report types, structures, and timelines, possibly referencing previously submitted instances. The system compiles the required information, reformats it according to the outlines provided, and initiates a submission-ready version to the validating agents. The report detector maintains a comprehensive list of report types together with their corresponding specifications. Validation agents perform various consistency checks at submission time to ensure these requirements are satisfied. Failure to perform such formatting checks may lead to major regulatory problems, such as the recent US\$ 2 trillion tax claim against the European Central Bank for breaching legal obligations without any intent to do so (Fleming, CH, 2022).

Timeliness and readiness controls are crucial factors for regulatory reports. Significant delays or late submissions can often trigger major fines. Although validation agents may enforce these rules, direct controls by the report detector are desirable. Natural Language Processing methods may assist in this regard by supporting an automatic description of the report's content and delay justification, also detecting possible deviations from the usual timing profile and alerting the internal clients.

4.3. Anti-Money Laundering and Fraud Detection Transaction monitoring combines real-time evaluation of payment flows with risk scoring determined through supervised learning models. A data lake collects structured transaction records and unstructured information from CRM systems and publicly available resources. Each transaction is characterized in terms of amount, frequency, purpose, counterparties, company profile, and risk status. Payments featuring high-risk factors — for example, large sums to or from offshore tax havens — may thereafter trigger special scrutiny and investigative case management, with a dedicated workflow for preliminary assessment and rule setting. Risk assessments incorporate external signals, like sanctions lists and negative news feeds, to enable automatic approval, rejection, or escalation for further inquiry.

Fraud detection applies similar principles as money laundering in terms of alert generation and investigation pathways, but is typically influenced by a stronger collaboration between business lines and line of defence functions. Use cases illustrate operational scenarios drawn from historical data, including purchase card usage, delegate spending, and trade finance outlier detection. Transaction monitoring models, deployed in production, have produced a satisfactory track record in terms of alert thresholds without generating an overload of cases.

4.4. Sanctions Screening and Export Control Compliance Sanctions Screening and export control compliance are critical components of any organization's control framework in trade finance and cross-border commerce. In such sensitive areas, where failure to comply with applicable regulations can lead to substantial reputational, operational, and financial damage, organizations typically invest significant sums in implementing the required compliance activities. Thereafter, such activity should be conducted on an exceptional basis only. Accordingly, these areas lend themselves particularly well to the application of agentic AI.

Sanctions screening involves checking screening names against actual transaction parties, which include clients, counterparties, consignees, consignors, settlement banks, and so forth. The list of sanctioned names, issued by local authorities in a national jurisdiction or by supranational authorities such as the European Union or the United Nations, is drawn from public and commercial sources. Updates are managed by the custodians and suppliers of the sanction lists. By definition, a clearmatch alerts the necessity for escalation and case management, whilst a no-match means that the screening requirement is satisfied. Inconsistency of result (i.e. a potential hit on one party name and a no-hit on an associated party name) should trigger abnormality detection and resolution. Screenings are also periodically refreshed for all names in the population, such that future higher-risk transactions do not go undetected in the interim.

5. Risk Management Frameworks for Self-Directed Systems

Governance and assurance mechanisms help to mitigate the risks of using self-directed models, enabling the world outside to better trust its outputs while also assuring the concerned institutions that they are operating within acceptable boundaries. Common areas of concern, and consequently, potential safety requirements are safety, explainability, human-in-the-loop decision-making, systematic verification throughout the development process, resilience, and contingency planning.

An external authority usually specifies safety objectives; they are thus context-specific. Financial institutions in the finance domain expect risk detection models to generate understandable warnings and explanations that users and regulators alike can evaluate. The classification systems that judge user

cardiovascular risks, or those related to the development of different diseases and their implications of granting or denying a loan, are some that model bias, fairness, and transparency.

The human-in-the-loop mode addresses the need for human oversight. Human decision-makers, crisis managers, or even flight controllers should always remain in a position to handle decisions that are either too critical or associated with a much higher level of risk than the decisions that are usually taken in an automated way. Another concern relates to the allocation of decision rights among the models and the users in the human-in-the-loop implementation mode, together with the definition of escalation protocols.

Verification, validation, and testing protocols guarantee the proper implementation and ongoing management of financial and risk models. Testing protocols can be organized in units, like the test of the engine of a plane in a testing facility, with well-specified expected results; integration testing, similar to cabin pressure testing with all systems off; system acceptance testing, like entire plane immersion testing; verification according to regulatory requirements; and, of course, ongoing monitoring. Resilience and contingency planning protocols establish, on the one hand, how to react when no model produces a reliable decision and what other elements the institution should put in place to quickly cope with possible failures in the decision-making process, and, on the other hand, how to cope with stress situations.

Table: Confusion-matrix table used to tune thresholds

	Predicted 1	Predicted 0
Actual 1	TP	FN
Actual 0	FP	TN

5.1. Safety, Explainability, and Auditability

Safety objectives encompass risk mitigation

during operation and accidents involving agentic AI. External safety considers third-party risks, while internal safety identifies faulty internal components and mechanisms to control damaging behavior. Agentic AI outputs must be explainable to assure decision-makers, action-takers, and affected stakeholders. The explanations should reveal sufficient information about the inputs and the state of the model producing the output, tailored to each audience, and may involve what-if analyses. Both the decisions and the underlying reasoning should be auditable to allow for retrospective accountability analysis by independent actors.

Several aspects related to auditability may be crucial for stakeholders affected by agentic AI decisions. Auditability demands that all operational aspects—from the data relied upon to the decisions taken and the monetary external effects—be available for analysis and assessment alongside a natural language description of the operation, including any processes that triggered it. The output must include details on data input sources, reliability and protection properties, data processing, and expected deviations from ordinary behavior; information about major changes since the last successful unit test; a history of related cases, if any; and appropriate escalation paths.

5.2. Human-in-the-Loop versus Fully Automated Decisions

Important distinctions arise within the definitions of agentic AI concerning human involvement in decision-making. Predominantly, controls exist at two points in the decision cycle. Firstly, within the intent

specification process, the governing human-assigned outcome can reside at the level of the corporate entity, the individual agent, or an intermediate position. Attributes of self-directed intelligence permit specification of intent at a higher or narrower granularity, depending on regulatory and business requirements. Confirmation of these aspects solidifies the argument for machine learning-based modelling, as these capabilities would be unattainable via traditional programming approaches.

Unlike intent determination, which can be fully automated, the point of decision execution remains a choice for the deployment architecture. Controls can be defined here to suit organizational preference and risk appetite. On the one hand, critical functions with high penalty or safety risk can be set to human-in-the-loop operation; decisions are made by the machine but require human approval to execute. Such controls can be tightened or relaxed, based on the ongoing output safety record of the decision pathway. Any undesirable event caused by a matching agent decision would automatically transition the decision pathway into a fully human-controlled mode. On the other hand, the system could execute fully automated decisions, even the most critical of functions with no human involvement. In this mode, the responsibility for negative consequences shifts from an individual to a collective of multiple independent agents, who acted according to defined risk controls. Although a negative event is still possible, it becomes far less likely—a trade-off that all firms must make based on their level of confidence in risk controls, which are ultimately grounded in historical data.



Fig 3: AI Agents With Human In The Loop

5.3. Verification, Validation, and Testing Protocols Verification, validation, and testing protocols encompass unit, integration, and acceptance testing. Individual components undergo unit testing to detect basic errors. Integration testing examines interoperability with other systems; the focus is on interface errors. For security-critical functions, regulatory verification guarantees adherence to prescribed standards. For example, Moneyval provides a list of criteria to assess states and jurisdictions for compliance with international AML/CFT standards. The entire system is subject to acceptance testing before being deployed into production and must pass regulatory scrutiny if required—such as for a medical product undergoing regulatory approval. During operation, monitoring assesses quality and timeliness.

In addition to standard testing, self-directed agents require additional protocol layers and specifications to maintain control over risk and safety. Safety verification assigns a safety objective to self-directed systems and confirms that they satisfy it in all scenarios. Safety considerations are particularly pertinent for physical agents (e.g., autonomous cars, unmanned drones, service robots). Systems in nonphysical domains (e.g. information processing) should still evaluate safety in terms of data leakage and breach, conflict with

regulators or other priority stakeholders, and dissemination of unsafe data or information. Monitoring must support early detection of violations and enable timely interruption of operations.

5.4. Resilience and Contingency Planning

Resilience and continuity planning for self-directed systems embrace diverse risks, specify planning objectives, and cover capability unavailability, extended unavailability, and geographic unavailability. Resilience describes the ability to endure, adapt to, and recover from disruption. As applied to self-directed systems, it encompasses recovery from incident-based disruptions (e.g., failures, security breaches) and evolving capability needs (e.g., training data for machine-learning components, training for human-in-the-loop decisions). For safety reasons, a system's conditions for safe operation should also be captured and monitored, ostensibly preventing dangerous operational states (e.g., generating gibberish, causing fatal crashes). Such states demand a distinct response. A recovery objective—recovering within a specified time period—is defined for an operational disruption.

Contingency plans define actions and responsibilities for when a system is unavailable or otherwise incapable of safe operation. Contingency planning considers three modes of non-availability: short-term loss of certain capability, extended unavailability (e.g., redesign), and localized loss of function in geofenced systems. In the first mode, a system is expected to be back in operation soon. In the second mode, either the whole system is down for a long time, or the downed component has a long lead time to repair. During a longer operation disruption, the lack of a critical system capability could expose the organization to safety, financial, or reputational risk.

Equation 4: Human-in-the-loop vs full automation (decision execution control)

A simple control policy (mathematical form):

Let p be model risk probability and let c be a criticality score (penalty/severity).

$$\text{RequireHumanApproval} = \begin{cases} 1 & \text{if } p \geq \tau \text{ or } c \geq c^* \\ 0 & \text{otherwise} \end{cases}$$

After an adverse event count E in a window:

$$\text{If } E \geq E^* \text{ then force human-only mode}$$

This encodes the paper's "transition into fully human-controlled mode" after undesirable events.

6. Ethical, Legal, and Social Implications

Self-directed intelligence has the potential to impact society on multiple levels. First, self-directed systems need to be designed, built, and operated such that they serve the public good and comply with legal requirements. Second, the impact of deploying such systems on human actors, their roles, and their skill requirements must be understood. Third, the implications of scaling to many self-directed systems, which can build other self-directed systems, need to be anticipated.

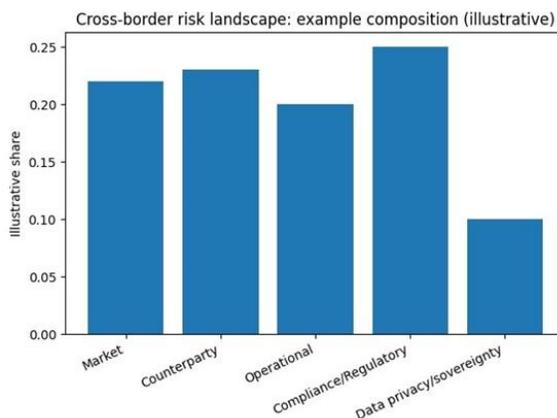
Inherent in the development and deployment of self-directed systems are questions of safety, accountability, bias, fairness, and societal impact. For every concern there are well-developed frameworks for mitigating risk, establishing accountability, addressing bias and fairness, and understanding wider societal impact. The AI principles of safety, accountability, and non-discrimination are the keystones of

responsible AI deployment. The risk management framework for self-directed intelligence proposed in this paper is, therefore, an elaboration of these principles.

Accountability encompasses instant responsibility for an action, liability for adverse consequences, and redress for affected parties. To ensure accountability, decision-making processes must be well-defined, traceable, transparent, and open to oversight. With current AI applications, there is little debate or uncertainty around accountability structures. The same cannot be said for self-directed systems, particularly those capable of authoring new and independent models with minimal human involvement. Delineating the boundaries of autonomy for interaction with human auditors, curators, and enforcers is vital to establishing accountability.

Ample research is dedicated to modelling and measuring bias and to establishing fairness criteria for AI applications. Continuous data labelling, especially for supervised models, raises the question of how representative the data remains. New data sources used to train and infer from models could introduce entirely new types of bias. An implicit assumption in the discussions of fairness in AI is that these models are still competitive and not bereft of drawbacks. However, these systems can be used for self-directed intelligence applications, where the new data is more representative of the trained population.

While the impact of AI systems has so far concentrated on narrow applications such as speech recognition and autonomous vehicles, the scaling of self-directed AI will influence how these systems interact in the future. It is not unreasonable to predict that self-directed bots will be able to create new self-directed bots of similar capability with minimal human help. These applications will drive structural market change, creating competitive advantage for early adopters.¹⁷⁸ Such a shift will necessitate the concurrent development of well-tested verification and validation mechanisms, as well as a robust set of AI principles. Agentic AI systems are often accused of being black boxes, but a significant number of these black-box aspects (such as data sovereignty or privacy algorithm use) are inherently a factor of the model and need not be hidden in discussions regarding the decision process.



6.1. Accountability and Liability

The assignment of responsibility for decisions made by agentic AI systems is an essential consideration, particularly in high-stakes environments where undesirable outcomes may draw blame. Since agentic AI can operate unattended over extended periods, establishing appropriate lines of responsibility is critical. In the case of self-directed systems, the rationale

and parameters guiding decision-making may be sufficiently clear to warrant accountability in the non-delegated sense. Moreover, the proximity of human oversight may further strengthen this connection. Nevertheless, it remains important to acknowledge that these systems ultimately possess the ability to act without direct human authorization. While for individual transactions the resulting decisions are hence delegated by design, day-to-day operations are subsequently overseen by human operators. During this period, agents indirectly assess the decisions being generated as well as the premise under which they are being made. Such oversight hence constitutes a human-in-the-loop arrangement, albeit in a passive role.

When supervision or oversight is insufficient, the absence of an intentional agent leads to questions of liability in the event of undesirable system behavior. Absent an agreement specifying otherwise, liability remains with the owner of the AI system, as in general consumer protection principles. Adopting a principle of risk allocation, organizations that derive economic benefit from AI systems should be legally responsible for undesirable consequences arising from their deployment and should also remunerate affected parties for material losses. Providing redress for detrimental impacts even when legal responsibility may not apply reinforces their accountability toward those affected, helping them protect their reputation and avoid social license issues when a wider stakeholder or community perspective is considered.

6.2. Bias, Fairness, and Transparency

Agentic AI systems, like any complex analytics tool, could introduce latent biases, resulting in unwarranted discrimination or lack of fairness in their operations. While bias in AI has received considerable attention in the literature, agentic AI is distinct. The major reason for this discrepancy is that, unlike supervised learning – an area where bias correction has been extensively developed – agentic AI relies fundamentally on unsupervised reinforcement learning. Thus, minimizing bias in agentic AI is arguably more similar to design for ethical or socially responsible behavior than to bias mitigation in supervised learning. Bias refers to the aggregation or error introduced into predictions on a population due to approximating a complicated function by a constrained model. Bias in unsupervised reinforcement learning is worse if the reward function has circularity (i.e., it rewards agents for achieving high rewards) or high fault toleration.

Despite these subtleties, it is advisable to scrutinize agentic AI models for possible bias. Bias indicates some uncontrolled influence over the decision process of any model, such that outcomes may be skewed in favor of certain classes of agents over others. Potential bias may arise in two ways. The first is direct bias, where, for example, the system is explicitly rewarding only the successful models. The second is indirect bias, where an external reward signal increases the likelihood of certain models being used by the system for completing tasks. Both scenarios must be monitored for any portion of a group being unduly managed in favor of others. Moreover, these controls need to be publicly stated with a description of why and how these operations are executed.

Once bias input has been adequately controlled, the final consideration is that of transparency. Transparency, distinct from explainability, refers to a more general view of the model workings, rather than a specific explanation for a single prediction, and thus requires an ex ante assessment of all potential inputs and decisions of the model being discussed.

6.3. Impact on Workforce and Firms

The widespread adoption of agentic models will reshape work practices and organizational structures within the finance industry. Executives and board members are responsible for determining the degree of decision-making automation. Instead of

implementing AI systems as mere assistants for human employees, companies can choose to automate entire decision-making processes and hand over control to the agents, provided that risk management procedures are in place and control processes performed regularly. Such utilization of AI agents might replace many human jobs and require different skills in the remaining positions in the company. Agentic models promise cost savings, higher efficiency, and knowledge that surpasses any individual. Nevertheless, new challenges also arise, as structured processes evolve into unstructured learning problems for the agentic AI models.

Many Agentic AI finance applications can be embodied into a different form for testing and consumption purposes, such as assuming the role of a fraudster or money laundering criminal. Agentic models might well follow the principles of adversarial training, where a model training to fulfill a useful role like detecting malware constantly learns from an attack agentic model with a limited rollout budget trying to evade detection.

The further use of agentic models ultimately depends on company and supplier choices. Agentic AI finance applications need to be tested and developed in the finance sector. Many research directions are conceivable. An emerging one is to diminish the training time of deep learning models.

7. Conclusions

Bringing together risk landscape analysis and compliance use cases elucidates the cross-border risk and compliance problem, while specifying agentic AI capabilities for compliance automation. Real-time monitoring, automated reporting, anti-money laundering and fraud detection, and sanctions screening are identified as concrete applications. A supporting risk management framework covers safety, explainability, and auditability criteria; human-in-the-loop and fully automated decision rights; verification, validation, and testing protocols; and resilience and contingency planning.

Tackling regulatory and market risk dimensions for financial institutions is a natural next step. Regulatory supervision is judiciously asserted as the appropriate source of safety, explainability, auditability, verification, and validation for self-directed systems. As organizations deploy such systems internally or build-as-a-service propositions, self-directed financial processes will complement more complex and high-value markets and counterparty risk models. Dialogue between financial-services incumbents, fintechs, and technology firms actively pursuing real-time monitoring, compliance automation, and risk mitigation affords industry actors a pioneering position in the space. Such distinct capabilities, pursued both independently and cooperatively, will navigate the previously cited challenges of regulatory heterogeneity, operational ineffectiveness, and market-fragility issues. A pragmatic response will progressively recruit testing and exploratory activity around market and counterparty risk.

7.1. Emerging Trends

Political dimensions will affect where AI is nurtured and which countries lead deployment. Supportive investment, talent retention, demand stimulation and responsibly calibrated regulation will favour growth. Supply chains will continue to be reshaped, and geopolitical alignments will transform. As independence is pursued, sectors such as semiconductors, batteries, pharmaceuticals and defence will be prioritised. National data are seen as a critical ingredient for economic and social prosperity. Countries will leverage their own data to enhance public services and analysis, while protecting data of economic and security interest. Inequality may be exacerbated both economically within countries and politically across countries.

The fusion of AI with finance will give rise to the substantive field of Self-Directed Intelligence and the third wave of AI, defined as custom-developed systems that can independently scope, learn, reason, decide and act in a sustained manner in specialised business domains and with organisational oversight. As human-directed automation transitions into more complex Self-Directed AI systems, some aspects of businesses may be subject to automated governance rather than just automated fulfillment. Self-Directed Intelligence will enable organisations to introduce an internal framework that embraces Agentic AI across the entire technology function by promoting systems that possess the capabilities required by cross-border risk and compliance processes.

8. References

- [1] Basel Committee on Banking Supervision. (2023). Principles for the effective management and supervision of climate-related financial risks. Bank for International Settlements.
- [2] Auer, R., Cornelli, G., & Frost, J. (2024). Artificial intelligence and machine learning in finance: Uses, risks, and supervisory considerations. BIS Working Papers.
- [3] Charoenwong, B., Kowaleski, Z. T., Kwan, A., & Sutherland, A. G. (2024). RegTech: Technology-driven compliance and its effects on profitability, operations, and market structure. *Journal of Financial Economics*, 154, 103792.
- [4] Charoen Wong, B., Kowaleski, Z. T., Kwan, A., & Sutherland, A. G. (2024). RegTech: Technology-driven compliance and its effects on profitability, operations, and market structure. *Journal of Financial Economics*, 154, 103792.
- [5] Davis, E. P. (2020). Financial stability, supervision, and agent-based systems. *Journal of Financial Regulation and Compliance*, 28(3), 321–337.
- [6] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT 2019* (pp. 4171–4186).
- [7] He, K., Gkioxari, G., Dollár, P., & Girshick, R. (2017). Mask R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision* (pp. 2961–2969).
- [8] Paleti, S. (2022). Fusion Bank: Integrating AI-Driven Financial Innovations with Risk-Aware Data Engineering in Modern Banking. *Decision Making*, 2326, 9865.
- [9] Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., & Narasimhan, K. R. (2024). SWE-bench: Can language models resolve real-world GitHub issues? In *Proceedings of the International Conference on Learning Representations*.
- [10] Kim, A. G. (2024). Financial statement analysis with large language models. *SSRN Electronic Journal*.
- [11] Lee, J., Lee, S., & Kim, J. (2024). A survey of large language models in finance (FinLLMs). *arXiv*. (arXiv:2402.02315)
- [12] Amistapuram, K. (2024). Privacy-Preserving Machine Learning Models for Sensitive Customer Data in Insurance Systems. *Educational Administration: Theory and Practice*. <https://doi.org/10.53555/kuey.v29i4.10965>

- [13] Liu, X., Yu, H., Tang, J., Gao, J., & others. (2024). AgentBench: Evaluating LLMs as agents. In Proceedings of the International Conference on Learning Representations.
- [14] Li, Y., Rao, J., Zhao, H., & others. (2023). Large language models in finance: A survey. arXiv. (arXiv:2311.10723)
- [15] Aitha, A. R. (2024). Generative AI-Powered Fraud Detection in Workers' Compensation: A DevOps-Based Multi-Cloud Architecture Leveraging, Deep Learning, and Explainable AI. Deep Learning, and Explainable AI (July 26, 2024).
- [16] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. Proceedings of NAACL-HLT, 4171–4186.
- [17] Nie, Y., Kong, Y., Dong, X., Mulvey, J. M., Poor, H. V., Wen, Q., & Zohren, S. (2024). A survey of large language models for financial applications: Progress, prospects and challenges. arXiv. (arXiv:2406.11903)
- [18] Reddy Segireddy, A. (2024). Federated Cloud Approaches for Multi-Regional Payment Messaging Systems. Turkish Journal of Computer and Mathematics Education (TURCOMAT), 15(2), 442–450. <https://doi.org/10.61841/turcomat.v15i2.15464>
- [19] Park, J. S., O'Brien, J. C., Cai, C. J., Ringel Morris, M., Liang, P., & Bernstein, M. S. (2023). Generative agents: Interactive simulacra of human behavior. In Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (pp. 1–22).
- [20] Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Zettlemoyer, L., Cancedda, N., & Scialom, T. (2023). Toolformer: Language models can teach themselves to use tools. In Advances in Neural Information Processing Systems (Vol. 36).
- [21] Shinn, N., Cassano, F., Berman, E., Gopinath, A., Narasimhan, K., & Yao, S. (2023). Reflexion: Language agents with verbal reinforcement learning. In Advances in Neural Information Processing Systems (Vol. 36).
- [22] Soria, J. J., & Romero, L. (2024). Machine learning models for money laundering detection: A review of approaches and challenges. In Proceedings of LACCEI 2024.
- [23] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In Advances in Neural Information Processing Systems (Vol. 30).
- [24] Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., Fan, L., & Anandkumar, A. (2023). Voyager: An open-ended embodied agent with large language models. arXiv. (arXiv:2305.16291)
- [25] Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q. V., & Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems (Vol. 35).
- [26] Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., & Cao, Y. (2023). ReAct: Synergizing reasoning and acting in language models. In Proceedings of the International Conference on Learning Representations.

- [27] Zhang, G., Li, X., & Zhou, Y. (2023). Machine learning approaches for constructing the national anti-money laundering index: Evidence from mutual evaluation reports. *Journal of International Financial Markets, Institutions and Money*, 82, 101671.
- [28] Zheng, L., Huang, X., & Liu, J. (2024). Deep learning for cross-border transaction anomaly detection: A temporal perspective. arXiv. (arXiv:2412.07027)
- [29] Zhou, W., Chen, Y., & Liu, S. (2024). Supply chain financial fraud detection based on graph representation learning. In *Proceedings of the International Conference on Knowledge Science, Engineering and Management* (pp. 112–126).