

Trustworthy Agentic Supply Chains: A Governance Framework for Digital Twin Orchestrated AI Decisioning Under Compliance, Auditability, and Data Sovereignty Constraints

Kannan Avalurpet Loganathan¹ and Arunraju Chinnaraju²

¹ Independent Researcher, California, USA.

² Doctorate in Business Administration, Westcliff University, USA.

DOI : <https://doi.org/10.51583/IJLTEMAS.2026.150100018>

Received: 14 January 2026; Accepted: 01 January 2026; Published: 23 January 2026

ABSTRACT

The rapid increase of Artificial Intelligence (AI) within Supply Chain Management (SCM), has transitioned SCM from using primarily Predictive Analytics & Decision Support, toward increased Autonomy of Decision Execution. However, although there are many examples of AI-driven SCM systems currently being used, they generally suffer from low levels of Trust, poor Governance structures, inadequate Auditability, and unresolved Data Sovereignty issues; all of which limit their potential deployment in High Consequence & Regulated Operational Environments. In an effort to address this important gap, this research introduces a comprehensive Governance First Framework for Trustworthy Agentic Supply Chains; where Autonomous AI Agents use Digital Twin Orchestrated Decision Intelligence to make decisions in accordance with explicit Compliance, Auditability, and Sovereignty Constraints. Agentic Supply Chains are defined as Socio Technical Systems, where Decision Authority is delegated to AI Agents that Continuously Sense, Simulate, Decide, and Act Across Dynamic Supply Networks. Digital Twins are redefined from Passive Visualization Tools to Active Orchestration Substrates that facilitate Real Time State Synchronization, Policy Execution, and Controlled Interaction Between Autonomous Agents and Enterprise Systems. On top of this base, the paper provides a Layered Reference Architecture for integrating Agentic Decision Intelligence, Bounded Autonomy, Governance by Design, and Human Oversight into a Unified Operational Model.

The Framework addresses Key Adoption Barriers via Explicit Mechanisms for Regulatory Alignment, Decision Traceability, Data Sovereignty Preservation, and Risk Containment. The Architectural Constructs provided include Agent Drift Detection, Rollback and Safe Recovery, Simulation Based Stress Testing, and Resilience under Adversarial and Extreme Disruption Scenarios. Through the embedding of Governance within the Decision Architecture, the proposed model allows Autonomous Supply Chain Systems to be Auditable, Compliant, and Strategically Controllable while Retaining Adaptive Intelligence. Additionally, beyond technical design, the paper outlines Evaluation Metrics, Organizational Integration Principles, Ethical Considerations, and Strategic Implications related to Delegating Decision Authority to Agentic AI Systems. Finally, the Study identifies a Forward-Looking Research Agenda addressing Multi-Agent Coordination, Cross-Enterprise Autonomy, and Next Generation Optimization Paradigms. Overall, this Work establishes Trustworthy Agentic Supply Chains as a Distinct and Necessary Evolution of Supply Chain Intelligence, providing a Reusable Reference Framework for Researchers, Practitioners, and Policymakers looking to operationalize Autonomous Decision Systems Responsibly and at Scale.

Keywords: Trustworthy Agentic Supply Chains, Agentic Artificial Intelligence, Digital Twin Orchestration, Autonomous Decision Intelligence, Governance-by-Design in AI Systems, AI Auditability and Decision Traceability, Data Sovereignty in Global Supply Networks, Bounded Autonomy and Human Oversight

INTRODUCTION TO TRUSTWORTHY AGENTIC SUPPLY CHAINS

Supply Chain Management has progressed from a focus on decision making authority through early supply chain systems focused on rule-based planning to predictive analytics and machine learning to current agentic artificial intelligence; each step represents an increasing level of sophistication in how decision-making authority is conceptualized and operationalized within supply chain production and distribution networks (Butner, 2010). Early supply chain systems used static heuristics, predefined thresholds and deterministic optimization routines (AlMulhim, 2021) based on stable demand patterns, predictable lead times and low levels of external volatility (AlMulhim, 2021). As global supply networks grew in terms of size and complexity, these assumptions became increasingly brittle (Butner, 2010). Predictive analytics and machine learning enhanced the ability to anticipate demand fluctuations, supply disruptions and capacity constraints through improved forecasting accuracy and situational awareness (AlMulhim, 2021). However, predictive analytics and machine learning generally have been applied as decision support systems, providing analytical insights to inform human judgment, as opposed to directly executing operational decisions (Butner, 2010). Therefore, the need for speed, adaptability and consistency in managing modern supply networks under continuous uncertainty remains constrained by human-mediated control structures (AlMulhim, 2021).

Decision support-oriented artificial intelligence becomes even less effective in environments with high-frequency disruptions, nonlinear interdependencies and cascading failure risk (Burgos & Ivanov, 2021). Predictive models can predict possible disruptions but they cannot inherently determine how competing objectives (such as cost, service, resilience, and compliance) will be prioritized in real-time (Burgos & Ivanov, 2021). Additionally, decision support systems suffer from latency because the insights generated by the system must be interpreted, escalated, and approved prior to taking action (Butner, 2010). The separation of prediction and execution creates structural delays that can limit the effectiveness of advanced analytics in environments with high levels of volatility (AlMulhim, 2021). Furthermore, decision support architectures do not scale as well as agentic architectures when the complexity of the network increases because human decision-makers are unable to effectively process and coordinate the volume of decisions that must be made across interconnected supply nodes (Butner, 2010). Together, these constraints highlight a fundamental gap between analytical intelligence and operational control in current supply chain systems (AlMulhim, 2021).

Agentic artificial intelligence presents a qualitative shift in this decision-making paradigm by introducing autonomous decision actors that are capable of perceiving system states reasoning about policy objectives and executing actions independent of human intervention (Jannelli et al., 2025). In agentic supply chains, decision authority is explicitly given to artificial agents that function within established areas of responsibility, such as inventory allocation, logistics routing, and supplier coordination (Mousa et al., 2024). Unlike decision support systems, these artificial agents do not only provide recommendations for action; they take action directly through enterprise systems based on policies that have been learned through experience and real-time feedback (Jannelli et al., 2025). This delegation of authority allows for ongoing adaptations to changing conditions and facilitates cross-layered coordination across various dimensions of operations (Mousa et al., 2024). Importantly, agentic systems represent intentionality in that they pursue defined goals over time instead of responding to singular events (Jannelli et al., 2025). Thus, artificial agents become persistent organizational actors that exist within the supply chain rather than as analytical tools that exist outside the realm of decision-execution (Butner, 2010).

The advent of agentic autonomy brings forth new organizational and control issues that go beyond technical optimization (Cheong, 2024). When artificial agents are granted the authority to make decisions autonomously, questions emerge concerning accountability, oversight, trustworthiness and compliance (Kacianka & Pretschner, 2021). Traditional governance structures in supply chains assume that decisions can be traced back to human managers whose actions can be audited, sanctioned or rectified through institutional processes (Kacianka & Pretschner, 2021). Agentic systems break this assumption by executing decisions continuously and at scale, often based upon complex learned representations that may not be easily understandable (Cheong, 2024). Without explicit governance structures, such autonomy can create opaque decision outcomes, violations of regulatory compliance, and unforeseen systemic behaviors (Phiri, 2025). These risks are exacerbated in global supply networks due to differences in data sovereignty, legal jurisdiction and ethical responsibility across regions (Hummel et al., 2021). Consequently, the primary issue is no longer whether artificial intelligence can optimize

supply chain decisions, but whether autonomous decision execution can be trusted, controllable and institutionally legitimate (Cheong, 2024).

This research is driven by the understanding that agentic supply chains must be designed using a governance-first approach as opposed to being viewed as an extension of predictive analytics or automation (Kacianka & Pretschner, 2021). Artificially intelligent decision-executing systems require that decision authority be bounded, monitored and auditable by design (Phiri, 2025). Digital Twin Orchestration represents a key mechanism for accomplishing this goal by allowing for real-time synchronization of state, simulated validation and controlled execution of agentic policies within a digital representation of the supply network (Busse et al., 2021). Through incorporating governance constraints, traceability mechanisms and sovereignty-aware data controls into the decision-execution architecture, agentic supply chains can balance autonomy with accountability (Abbas et al., 2024). The unique contribution of this study lies in reframing autonomy as an organizational control issue, rather than a strictly computational optimization issue, and in developing a comprehensive framework that allows artificial agents to serve as accountable decision-executors within complex regulated and high-consequence supply chain environments (Sani et al., 2024).

Conceptual Foundations of Agentic AI in Supply Chains

Artificial intelligence that acts as an entity in decision-making processes for extended periods of time to make decisions, instead of being used as analytical tools (Jannelli et al., 2025) is referred to as "agentic artificial intelligence." In supply chain management, the defining characteristic is not just predictive ability; however, it is also decision-making in the long-term sense that the decision-maker continually assesses the current operating condition, forms an internal representation of the current condition, selects an appropriate course of action and assesses the outcome of that course of action relative to specific goals (Xu et al., 2024). The reason why the decision-maker's continuous assessment of the environment is so important is that the supply chain is not an optimization problem that remains constant; it is an evolving network of relationships between all participants in the supply chain, which is influenced by many factors such as feedback, delay, constraints and disruptions (Butner, 2010). Therefore, agentic artificial intelligence is better thought of as an operational participant in the control structure of the supply network, able to act and coordinate its own responses to unknowns and remain compliant with institutional requirements (Dominguez & Cannella, 2020).

The following properties differentiate agentic AI from standard supply chain analytics (Rolf et al., 2023): Goal-directedness enables decision-making to be made in terms of multiple horizons, in that the decision maker pursues objectives such as maintaining service levels, controlling costs, increasing resilience, and achieving regulatory compliance as separate yet interdependent priorities across time (Xu et al., 2024). Situational awareness enables continuous perception, integrating real-time data from various sources, including demand signals, supplier availability, transportation capacity, operational status, contracts, and policies (Terrada et al., 2020). Adaptability enables the decision maker to learn from its previous experiences and adjust its decision strategy when the environment changes (e.g., due to changes in seasonal patterns in demand, supplier reliability, port congestion, or regulatory changes) (Zhang et al., 2024). Boundedness limits the extent to which the decision maker can act, acknowledging that real-world supply chains are subject to limitations related to safety, contractually agreed upon requirements and regulatory requirements that may limit the short-term optimal use of resources (Papagiannidis et al., 2025).

In order to achieve conceptual clarity, it is essential to distinguish between automation, autonomy, and agency (Butner, 2010), since each implies a different level of authority and accountability for decision-making. An automated system executes predetermined procedures and is usually depicted as conditional logic or as a deterministic workflow that triggers actions once certain conditions have been met (Dominguez & Cannella, 2020). An automated system does not have an internally defined objective function that is optimized over time and does not adjust the logic used to determine action based on past experiences (Terrada et al., 2020). Autonomy extends beyond automation by allowing the system to decide which action to take from a set of predefined choices, and often uses optimization or learning-based prediction to identify the best possible action given a representation of the state (Rolf et al., 2023). Agency extends beyond autonomy by incorporating the authority to create new actions, pursue goals across time, adapt decision strategies based on feedback received from the

environment, and coordinate decisions across interconnected domains (Xu et al., 2024). In a supply chain setting, agency implies that the decision-making process is not limited to selecting the most desirable action given a locally defined decision-making function, but rather determines the trajectory of operations through continuous decision loops (Jannelli et al., 2025).

Delegating decision authority is the primary method through which agency becomes operational in supply chains (Xu et al., 2024). Delegation involves the deliberate transfer of decision authority to an artificial agent under clearly defined conditions and guidelines regarding when to escalate a decision to a human decision-maker, the scope of decision authority and how accountability will be established (Papagiannidis et al., 2025). The delegation concept is critical because agentic supply chains do not diminish human responsibility, but rather distribute decision-making responsibilities across temporal and cognitive scales (Novelli et al., 2024). High-frequency operational decisions (such as inventory rebalancing, carrier selection, replenishment timing, and exception handling) frequently require rapidity and consistency that exceeds the capabilities of humans (Butner, 2010). Through delegation, those types of decisions can be performed by an agent, whereas humans maintain authority over strategic parameters, constraint definition, risk tolerance and override authority (Raji et al., 2020). As a result, the effectiveness of an agentic system is dependent on the institutional mechanisms employed for delegation as much as the predictive accuracy and optimization performance of the underlying algorithms (Papagiannidis et al., 2025).

An agentic decision-maker's behavior in supply chains should be viewed as an adaptive controller in a coupled dynamic system (Garvey et al., 2015). A supply chain state represents the changing configurations of the following variables: inventory positions, capacity assignments, lead times, transportation network conditions, backlog orders, supplier performance and policy constraints (Rolf et al., 2023). When an agent takes an action, the state changes; the changed state subsequently alters the probability distribution of potential future occurrences (Oroojlooyjadid et al., 2022). Due to this feedback relationship, there exists path dependency; i.e., early decisions can alter the feasible options for subsequent decisions, possibly generating additional risks by creating cascading effects (e.g., stockout, expedited shipping, overtime labor, penalty clauses) (Chaharsooghi et al., 2008). Therefore, agentic decision-making is inherently sequential and requires thinking about the future consequences of an action and not solely selecting the optimal action at a point in time (Kim et al., 2024).

An expected utility framework is one way to concisely represent the decision logic of an agentic decision-maker (Oroojlooyjadid et al., 2022). Expected utility calculates the sum of the expected utility of a particular action, expressed in terms of the possible future states of the world, each weighted by its associated probability of occurrence given that the action has been taken (Chaharsooghi et al., 2008). Mathematically, the framework can be expressed as follows (Rolf et al., 2023).

$$\mathbb{E}[U] = \sum_{s \in \mathcal{S}} P(s | a) R(s, a)$$

The above expression can be viewed as follows; the Expected Utility [E[U]] denotes the expected utility for selecting action 'a', S denotes the set of all relevant system states, P(s|a) denotes the probability of achieving state 's' given action 'a' and R(s,a) denotes the rewards or value that results from the state-action pair (Oroojlooyjadid et al., 2022). As stated by Kim et al. (2024), in supply chains, the reward term is generally a composite of several metrics: cost, service level, timeliness, compliance adherence and risk exposure (Papagiannidis et al., 2025). Therefore, the probability term is similarly complex due to the variety of sources of uncertainty that exist in supply chains, including demand variability, supplier disruption, transportation delay and policy change (Burgos & Ivanov, 2021). The conceptual value of the equation is that it shows that agent-based decision making is not a single-step optimization, but rather a probabilistic assessment of the implications of taking an action relative to an organization's goals (Xu et al., 2024).

When contrasting the use of inference versus control, the distinction between Agent Based Decision Making and traditional predictive analytics will become even clearer (Rolf et al., 2023). Traditional predictive analytics primarily focuses on forecasting future events and quantifies future events including lead time, demand,

equipment failures, delay risks etc. (Butner, 2010). While predictive analytics can produce extremely high-quality predictions, the inability of organizations to react in a timely manner can render these predictions useless (Dominguez & Cannella, 2020). On the other hand, agent-based decision-making systems incorporate the control process by integrating both the inference process and the selection of an action along with the execution of the action (Jannelli et al., 2025). The agent-based decision-making system utilizes predictive signals as input to a policy that determines what action to take, when to take it and how to coordinate among various nodes (Xu et al., 2024). By incorporating the control function into the decision-making process, agent-based decision-making systems reduce the decision-making cycle time and increase the predictability of the organization's responses, allowing the organization to adapt to changing conditions more quickly than would be possible using human decision-making alone (Terrada et al., 2020).

Agent-based decision-making systems differ from traditional optimization methods in the way they define and implement the objective functions and constraints (Papagiannidis et al., 2025). Optimization problems traditionally rely upon a fixed objective function, well defined constraints, and full knowledge of the situation at the start of the planning period (Butner, 2010). However, real-world supply chains consistently violate these assumptions due to limited visibility, delayed access to information, and continually changing constraints (including changes to contracts, changes to compliance obligations, and capacity shocks) (Garvey et al., 2015). Agent-based decision-making systems can operate effectively under partially observable conditions by developing and maintaining internal beliefs about the world, updating those beliefs as additional information arrives and selecting actions that are robust against uncertainty (Rolf et al., 2023). Constraints are not viewed as static inputs but rather as parameters that evolve over time and constrain the range of available actions (Novelli et al., 2024). This approach views the governance and definition of constraints as primary theoretical constructs as opposed to implementation details (Raji et al., 2020).

An important aspect of the conceptual framework is the temporal nature of decision making in supply chains (Butner, 2010). A number of key supply chain outcomes arise from the accumulation of effects over time, such as progressive inventory depletion, compounding delays, and demand amplification (Burgos & Ivanov, 2021). To accommodate this temporal structure, agent-based decision-making systems optimize decisions as a sequence of decisions rather than as a series of discrete decisions (Oroojlooyjadid et al., 2022). This means that agent-based decision-making systems need to have the capability to assess the long-term impacts of the decisions made today, such as how today's expediting decisions impact tomorrow's cost base, supplier behavior and service volatility (Kim et al., 2024). Thus, while agent-based decision-making systems represent a significant improvement over traditional decision-support systems in terms of speed and responsiveness (Dominguez & Cannella, 2020), they represent a fundamentally different decision-making paradigm, in that continuous policy-driven execution constitutes the primary mechanism for maintaining the stability and performance of the supply network under uncertainty (Xu et al., 2024).

Another area in which agent-based decision-making systems provide a conceptual differentiator is coordination (Xu et al., 2024). Supply chains consist of numerous functional areas, including procurement, production, warehouse management, transportation and fulfillment, each with its own specific constraints and performance measures (Lee et al., 2008). Traditional analytics generally provide optimizations within individual silos or provide coordination via periodic planning cycles (Butner, 2010). Agent-based decision-making systems provide coordination across functional areas by viewing the supply chain as a shared environment in which multiple decision entities act (Terrada et al., 2020). Coordination can occur through shared state representations, hierarchical decision structures, or negotiated action protocols (Dominguez & Cannella, 2020). The conceptual implication of this is that performance is not dependent on the degree of "intelligence" of a particular model but is instead dependent on the structure of interactions among the agents, the degree of alignment of the objectives of the agents and the governance rules that regulate the potential for adverse competition or oscillatory behavior (Zhang et al., 2024).

Once decision authority has been delegated to the agents, it is natural to consider questions related to trustworthiness (Novelli et al., 2024). Trustworthiness in this case refers to the institutional property that emerges from the alignment of agent behavior with the intended purpose of the organization, the applicable regulatory requirements and the accountability mechanisms that govern the behavior of the agents (Papagiannidis et al.,

2025). Agents may be capable of producing very accurate predictions, producing high levels of cost savings and still be untrustworthy if their decisions are opaque, unverifiable, or non-compliant with applicable regulations (Raji et al., 2020). Trustworthiness thus becomes a characteristic of the entire socio-technical system, including the specification of the constraints, the specification of the escalation pathways, the logging of decision provenance and the capability to recreate and evaluate the rationale for the decisions made during execution (Novelli et al., 2024). Thus, this conceptual shift enables researchers to broaden the scope of supply chain research from a narrow focus on the achievement of performance metrics to a broader focus on governable autonomy (Raji et al., 2020).

In total, these conceptual foundations establish agent-based decision-making in supply chains as a distinct research area that goes beyond automation and predictive analytics (Jannelli et al., 2025). The defining difference is the delegation of decision authority to continuous policy driven execution, which transforms supply chain control into an adaptive system problem that is inextricably linked with governance, accountability, and legitimacy (Xu et al., 2024). By establishing clear definitions and distinguishing between automation, autonomy, and agency, the conceptual framework establishes the necessary foundation for serious discussions regarding how delegated decision-making agents behave in dynamic supply networks, how the objectives of these agents should be constructed, and how trustworthiness can be established as a measurable characteristic of the entire system as opposed to merely a desired managerial outcome (Papagiannidis et al., 2025).

Digital Twin Orchestration as the Decision Execution Substrate

Descriptive digital twins provide situational awareness through visualizations of supply chain activities and asset flows; however, they remain passive artifacts and do not support decision-making authority (Kritzinger et al., 2018; Negri et al., 2017). Since decision authority was located outside the digital twin, with either a planner or downstream analytical tool providing recommendations for action, the descriptive digital twin functioned as an observational artifact, and not an active control mechanism (Negri et al., 2017) – its primary limitation being insufficient for decision support during periods of extreme supply chain volatility (Burgos & Ivanov, 2021).

A key conceptual distinction is that descriptive digital twins differ from operational digital twins (Fuller et al., 2020). Descriptive digital twins are focused on representing the state of the system for purposes of interpretation and analysis, whereas operational digital twins are focused on influencing the state of the system (Tao et al., 2018). An operational digital twin is capable of internally updating its own state in real-time, and in doing so, is tightly-coupled to interfaces for validating decisions and executing actions (Wang et al., 2022). This tight-coupling between the digital twin and AI-driven decision logic and physical supply chain operations, enables the digital twin to function as an intermediary between those two spaces (Ivanov & Dolgui, 2021). Therefore, in the context of agentic supply chains, the operational digital twin represents the substrate upon which autonomy is exercised in a controlled and auditable manner (Raji et al., 2020). Moreover, the transformation from a digital twin functioning as a "mirror" of reality, to one that functions as an "execution environment," fundamentally alters how decisions affect the physical world (Grieves & Vickers, 2017).

The ability of an operational digital twin to function as a decision-execution substrate is predicated upon the existence of real-time state-synchronization capabilities (Fuller et al., 2020). The state of a supply chain is constantly evolving due to changing order demand, shipment movement, consumption of capacity, and disruption events (Ivanov & Dolgui, 2021). As such, the state of the supply chain is reflected through a myriad of disparate data streams that originate from a variety of sources, including enterprise systems, sensors, logistics platforms, and partner networks (Lee et al., 2015). The digital twin consolidates these disparate data streams into a coherent state representation that accurately captures the current configuration of the supply network (Wang et al., 2022). For state-synchronization to occur, it is not sufficient to merely ensure temporal alignment of the state representations of the digital twin and the physical supply chain. Additionally, there must exist semantic consistency across decision contexts regarding how inventory quantities, lead-times, and capacity constraints are interpreted (Tao et al., 2019). If such semantic consistency does not exist, autonomous agents will make decisions based on inaccurate or conflicting information, which will undermine both the performance and trustworthiness of the decision-making process (Novelli et al., 2024).

State-synchronization can be formally expressed as a state-update function in which the internal state of the digital twin is updated based upon the receipt of new observations (Fuller et al., 2020). If we denote the internal state of the digital twin at time t as x_t , and if we denote the set of new observations received from the physical system at time t as o_t , then the state update can be represented as follows:

$$x_t = g(x_{t-1}, o_t)$$

Where g denotes the state-reconciliation function that combines new observations with prior state estimates to produce a consistent state representation (Tao et al., 2019). The state-reconciliation function is critical to the development of accurate state representations, and the accuracy of the state representation has significant implications for the reliability of downstream agentic decision-making (Ivanov & Dolgui, 2021).

Additionally, digital twins enable event-driven simulation that serves as the second key capability that enables digital twins to evolve from descriptive artifacts to operational control substrates (Kritzinger et al., 2018). Supply chain actions result in cascading effects that occur over time and space (Ivanov & Dolgui, 2021). A decision to expedite shipments, alter production schedules, or reallocate inventory will impact downstream availability, upstream replenishment signals, transportation utilization, and contractual performance (Burgos & Ivanov, 2021). Event-driven simulation enables the digital twin to model how state transitions occur in response to actions, thereby enabling the digital twin to simulate the effects of actions (Tao et al., 2018). Consequently, event-driven simulation enables the digital twin to react proactively to potential disruptions, and not reactively once the disruptions have occurred (Fuller et al., 2020). This proactive capability is particularly important in agentic supply chains where decisions are made and executed at machine-speed (Wang et al., 2022).

Simulation within the digital twin supports decision evaluation by estimating future trajectories based on alternative actions (Tao et al., 2018). We can represent the projected future state as follows, where a_t denotes a candidate action proposed by an agent at time t , and f represents the system-dynamics encoded within the digital twin:

$$x_{t+1} = f(x_t, a_t)$$

This representation highlights that actions are evaluated within the digital twin prior to their impact on the physical system (Fuller et al., 2020). The function f incorporates physical constraints, policy rules, and environmental responses to determine whether an action violates feasibility or governance constraints, and if so, what modifications should be made to the action (Reichert & Weber, 2012). Therefore, simulation can serve as a gatekeeping mechanism that ensures autonomous decision-making aligns with organizational intent and operational realities (Novelli et al., 2024).

Furthermore, beyond single-step projections, event-driven simulation enables the estimation of multi-step trajectories across time horizons (Kritzinger et al., 2018). This capability enables the digital twin to estimate cumulative effects resulting from sequential actions, such as inventory depletion, service degradation, and/or cost escalation (Ivanov & Dolgui, 2021). By simulating these trajectories, agents can select actions that maximize long-term outcomes, rather than responding myopically to short-term signals (Wang et al., 2022). Temporal reasoning, facilitated by the digital twin-mediated execution process, distinguishes digital twin-mediated execution from rule-based automation and enhances stability in highly dynamic environments (Tao et al., 2019). Furthermore, temporal reasoning reduces the likelihood of oscillatory behavior, whereby rapid reactions exacerbate volatility, rather than dampen it (Burgos & Ivanov, 2021).

Additionally, digital twins function as mediators between decision-intelligence generated by autonomous-agents and enterprise-execution systems (Lee et al., 2015). In a governed architecture, agents do not invoke execution commands on enterprise-platforms directly (Raji et al., 2020). Instead, proposed actions are instantiated within the digital twin, and are subsequently validated against constraints, including capacity limits, contractual commitments, regulatory requirements, and risk thresholds (Reichert & Weber, 2012). Actions that meet these constraints are converted into executable commands (Wang et al., 2022). The mediation-layer established by the digital twin separates learning-processes from execution-interfaces, thereby preventing direct instability of operations resulting from policy-exploration and adaptation (Novelli et al., 2024). Thus, the digital twin

represents a separation between decision-reasoning and physical actuation, which is necessary for establishing trustworthiness in autonomous-decision-making (Kacianka & Pretschner, 2021).

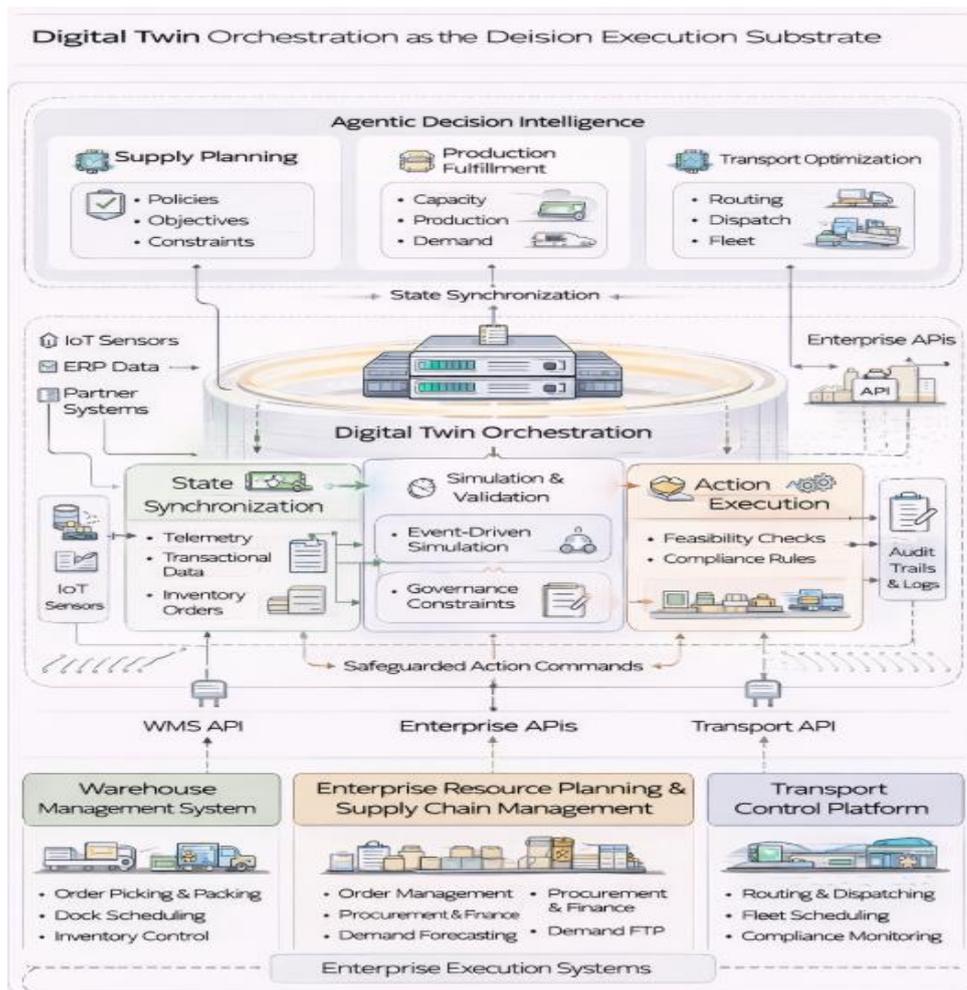
Furthermore, the mediation-layer provided by the digital twin facilitates coordination among multiple autonomous-agents operating in different functional-domains (Monostori et al., 2016). In agentic-supply-chains, procurement, production, logistics, and fulfillment agents may pursue objectives that interact through shared resources and constraints (Ivanov & Dolgui, 2021). The digital twin represents a shared-state-representation and orchestration-environment where competing actions can be evaluated collectively (Wang et al., 2022). Through orchestration, the digital twin resolves conflicts, prioritizes actions, and enforces system-wide policies that supersede individual-agent-objectives (Reichert & Weber, 2012). Hence, the digital twin transforms the supply-chain into a coherent multi-agent-system, rather than a collection of independent decision-silos (Van der Aalst, 2016).

Finally, from a governance-perspective, the orchestration-capabilities of the digital twin facilitate accountability and traceability by integrating decision-evaluation into the execution-pathway (Raji et al., 2020). With each decision, a sequence of state-assessments, simulations, and validations occurs, and these events can be logged as part of an execution-log (Van der Aalst, 2016). The execution-log enables post-execution reconstruction of decision-rationale, including the state-conditions considered, actions evaluated, and constraints enforced (Kacianka & Pretschner, 2021). This traceability is essential for demonstrating compliance with regulatory and contractual obligations in autonomous-environments (Novelli et al., 2024). Descriptive-twins do not possess this capability, as they observe outcomes without capturing the internal decision-process that produced those outcomes (Negri et al., 2017).

The orchestration-role of digital twins also transforms the temporal-structure of supply-chain-control (Fuller et al., 2020). Traditional-planning-processes operate on periodic-cycles, resulting in latency between observing the state-of-the-system and taking-action (Ivanov & Dolgui, 2021). However, digital-twin-mediated-execution enables continuous-decision-loops where state-changes immediately initiate evaluation and response (Wang et al., 2022). This temporal-compression improves responsiveness while maintaining stability through simulation-based-validation (Tao et al., 2019). The digital-twin functions as a stabilizing-buffer that reconciles the speed of autonomous-decision-making with the caution required for high-consequence-operations (Kacianka & Pretschner, 2021).

Ultimately, the digital twin becomes the locus where strategic-objectives, operational-constraints, and regulatory-requirements converge into executable-control-logic (Raji et al., 2020). Organizations can delegate operational-autonomy, while retaining accountability, through the centralized mediation of decision-making within the digital-twin (Novelli et al., 2024). Thus, digital-twin orchestration represents foundational-infrastructure for agentic-supply-chains, rather than an auxiliary-visualization-technology (Grieves & Vickers, 2017). It enables responsible-autonomy-at-scale, while preserving-trust, control, and institutional-legitimacy (Ivanov & Dolgui, 2021).

Figure 1: Digital twin supply chain architecture



This architectural model shows how an operational digital twin can be used as a central control and execution platform that serves as a bridge between decision intelligence (agentic) and enterprise execution systems, to transform a digital twin from a passive reflection tool into an active control mechanism. The upper layer consists of multiple, independent, decision-making domains such as supply planning, production fulfillment, and transport optimization, all producing policy-driven action recommendations based upon objective functions, constraint functions and learned strategies. None of these agents will act upon physical systems. Their proposed actions and state queries will flow down to the digital twin's orchestration layer through continuous synchronization of states. Within the digital twin, the various disparate real-time data streams from IoT sensors, telematics, ERP transactions, partner systems and inventory records, are merged into a single, semantically-aligned representation of the supply chain state. Therefore, when the agents make decisions, they are always grounded within a coherent and up-to-date view of the entire system. Using an event-based simulation and validation module within the digital twin, proposed actions are evaluated based upon their downstream effects propagated through encoded system dynamics, capacity constraints, governing rules and compliance policies, allowing for predictive analysis of potential congestions, delays, inventory imbalances and/or policy infractions prior to the implementation of those actions. The simulation-based validation step ensures that the proposed actions pass both feasibility tests and risk thresholds to ensure that only those actions that meet both the operational, contractual and regulatory requirements are approved. Approved actions are subsequently transformed into safeguarded, deterministic control commands and sent via standard Enterprise Application Programming Interfaces (APIs) to the appropriate execution platforms including Warehouse Management Systems, Supply Chain Management Systems, Enterprise Resource Planning Systems and Transport Control Platforms. These execution systems remain responsible for executing physical activities, i.e., ordering, dock scheduling, purchasing, routing, and compliance enforcement, while the digital twin remains responsible for maintaining transactional consistency and coordinating the different domain-based execution systems. Throughout this process, the digital twin tracks changes in states, the results of the validation processes and the

results of the execution decisions made into audit trails and logs and therefore provides transparency and accountability to the autonomous decisions made. As such, the digital twin represents a closed-loop control system where the actual-world execution continually updates the twin's state, the agents update their policies based upon the synchronized feedback from the twin and the twin acts as a buffer to slow the rate at which the agents learn compared to the rate at which they execute through simulation-based caution. Ultimately, the architecture described above represents a hierarchical but tightly-coupled control architecture that allows for real-time, decentralized, agential decision-making while maintaining corporate-wide responsibility/accountability/safety/alignment across complex supply chains.

Reference Architecture for Trustworthy Agentic Supply Chains

The reference architecture for a trustworthy Agentic Supply Chain (ASC) creates a structural framework that integrates the attributes of autonomy, trust and governance as a unified system property (Raji et al., 2020). Instead of viewing autonomy as an emergent product of advanced analytics, the reference architecture defines how decision authority is allocated, enforced, constrained, and monitored throughout the supply chain (Papagiannidis et al., 2024). Due to the inability of post-execution monitoring and organizational policy to effectively govern autonomic behavior, it is crucial to provide a structured representation of trustworthy autonomy in the connection between data perception, decision reasoning, action execution and oversight (Mitchell et al., 2019).

The Reference Architecture for ASCs is composed of multiple layers. The Data Ingestion and Contextualization Layer serve as the foundational layer of the reference architecture (Batini et al., 2009). Supply Chains generate large amounts of unstructured and heterogeneous signal types (transactional record; sensor telemetry; logistics event; communication from partners; regulatory data), (Koot et al., 2021). However, raw ingestion of these signals is insufficient for enabling agents to make autonomous decisions, as autonomous agents require semantically consistent and contextually relevant representations of the current state of the system (Simmhan et al., 2005). The contextualization of diverse signals enables the creation of a consistent view of the state of the system through the creation of a standardized and interpretable state representation (e.g., effective inventory availability; capacity commitments; lead time distributions; compliance constraints), (van der Aalst, 2016). Additionally, the resolution of temporal alignment enables a common view of the system at each point in time (Abideen et al., 2021). Therefore, the Data Layer enables agents to reason about the system with a level of situational awareness similar to that possessed by the organization and not just individual metrics (Hummel et al., 2021).

In addition to integrating data, the Data Layer also performs normalization, validation, and uncertainty representation (Batini et al., 2009). Supply Chain data is often delayed, incomplete, and/or noisy (due to delays in reporting, differences in data reporting practices between partners, and operational disruptions), (Chandola et al., 2009). Therefore, the Data Layer is also responsible for estimating confidence in data representations and detecting anomalies in order to inform downstream decision-making logic (Chandola et al., 2009). The architecture does not mask uncertainty, but instead presents it as an explicit component of state representations (Zhou et al., 2025). This enables agents to adjust their decision-making aggressiveness in response to lower levels of confidence in state representations (Xia et al., 2020). Therefore, the Data Layer is not only a conduit for information flow, but also a control mechanism that influences both the quality and reliability of autonomous decision-making (Simmhan et al., 2005).

The Decision Intelligence Layer transforms contextualized state representations into executable policies (Oroojlooyjadid et al., 2022). The Decision Intelligence Layer comprises of various components (forecasting, optimization, learning, and reasoning) that compare alternative courses of action against organizational objectives (Kim et al., 2024). Unlike traditional analytics, the Decision Intelligence Layer does not simply create recommendations for course of action, but creates executable policies that can be autonomously acted upon (Oroojlooyjadid et al., 2022). Organizational objectives are represented as multi-dimensional value functions that balance cost, service, resilience, compliance and risk exposure (Xia et al., 2020). Constraints representing governance and regulatory requirements are integrated into policy evaluation (Feinberg & Schwartz, 1995).

Therefore, the architecture ensures that intelligence and governance are not separate concerns, but co-determinants of decision outcomes (Raji et al., 2020).

Unlike traditional analytics, the Decision Intelligence Layer in Agentic Supply Chains operates continuously rather than episodically (Oroojlooyjadid et al., 2022). Policies are revised in real-time as new information becomes available, and as outcomes indicate changes in the dynamics of the system (Kim et al., 2024). Continuous adaptation requires mechanisms to stabilize the decision-making process to avoid oscillating behavior or reacting too aggressively to transitory signals (Xia et al., 2020). Stability is ensured through policy smoothing, horizon-based evaluation, and constraint-aware optimization (Feinberg & Schwartz, 1995). The architecture ensures that learning processes remain aligned with organizational intent by separating policy generation from execution and validating proposed actions downstream (Laato et al., 2022). In this manner, the Decision Intelligence Layer produces intent, rather than action, and retains execution authority for downstream layers (Papagiannidis et al., 2024).

The Agent Execution and Coordination Layer is the location of autonomy within the architecture (Kim et al., 2024). Agents operating within this layer are responsible for executing decisions within delegated domains (inventory allocation; transportation routing; supplier selection, etc.), (Oroojlooyjadid et al., 2022). Each agent has a defined action space, authority boundary and escalation protocol (Raji et al., 2020). The coordination of agents is necessary, since supply chain decisions interact with one another through shared resources and constraints (Kim et al., 2024). The architecture facilitates coordination through shared state representations, hierarchical control relationships and arbitration mechanisms that resolve conflicts among competing actions (Monostori et al., 2016). This coordination ensures that local optimizations do not undermine overall system performance (Koot et al., 2021).

Agent execution within the architecture adheres to the principles of Bounded Autonomy (Feinberg & Schwartz, 1995). Agents have the authority to execute actions independently within predetermined boundaries, but are subject to constraint enforcement and oversight (Haskell & Jain, 2013). Proposed actions of agents are not executed directly upon physical systems. Rather, proposed actions are forwarded through validation mechanisms to determine whether actions meet feasibility, compliance and policy criteria (Raji et al., 2020). This separation enables agents to execute actions rapidly, while maintaining safeguards to protect the organization from unintended consequences (Papagiannidis et al., 2024). The Execution Layer, therefore, realizes autonomy without relinquishing control (Laato et al., 2022).

The formal expression of agent execution within the architecture can be formulated as a constrained action selection process (Feinberg & Schwartz, 1995). Let the symbol a^* represent the action selected for execution and let A represent the set of all feasible actions given the current state x and governing constraints C . The execution rule can be expressed as

$$a^* = \arg \max_{a \in A(x,C)} V(x, a)$$

where V represents the value function generated by the Decision Intelligence Layer (Xia et al., 2020). The above formulation emphasizes that agent execution is conditioned on the feasibility and governing constraints of the system (Haskell & Jain, 2013). The architecture therefore ensures that agents exercise autonomy within an explicitly bounded decision space (Feinberg & Schwartz, 1995).

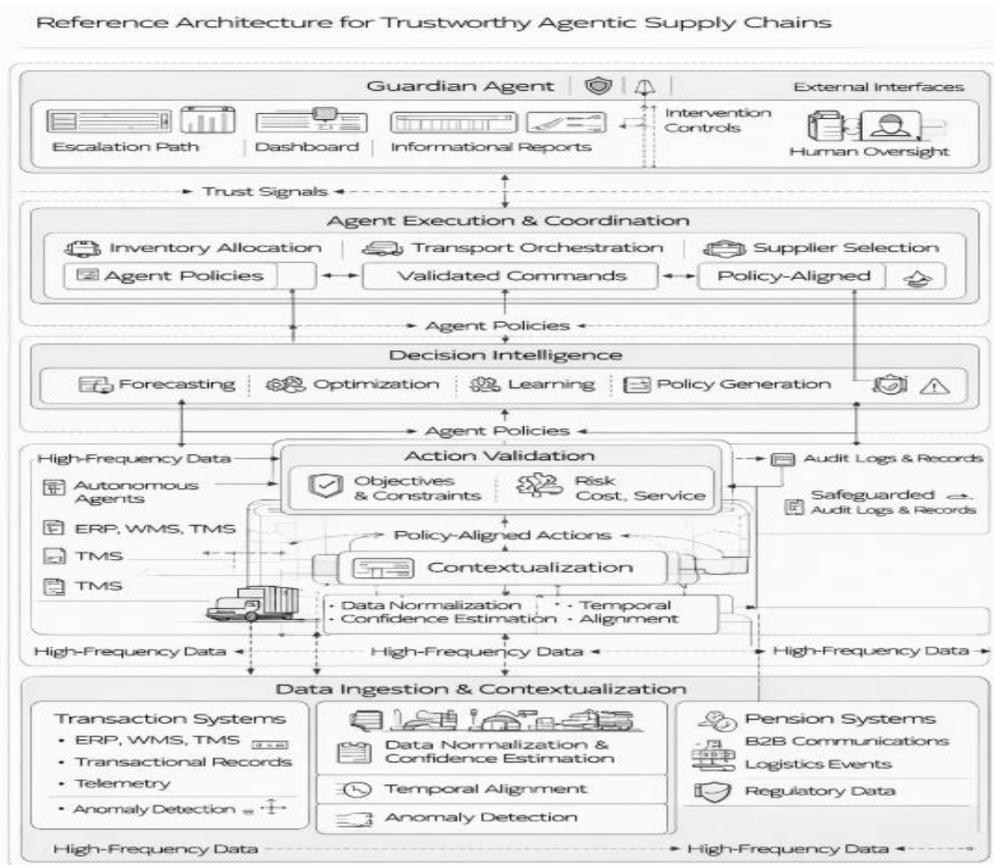
The Governance and Oversight Layer form the Normative Backbone of the Reference Architecture (Papagiannidis et al., 2024). Governance is viewed as an integral component of the execution pathway and is not solely viewed as an audit function (Raji et al., 2020). The Governance Layer defines policies; defines constraints; defines escalation thresholds; and defines the mechanisms of accountability that guide both the evaluation of and enactment of decisions (Laato et al., 2022). Governance Rules may define Regulatory Requirements; Contractual Obligations; Ethical Principles; and Organizational Risk Tolerances (Hummel et al., 2021). By embedding such rules in the architecture, the system is able to ensure that autonomous decisions continue to be institutionally legitimate, even as they evolve over time in response to changing conditions (Mitchell et al., 2019).

Oversight Mechanisms within the Governance Layer include Decision Logging; Traceability; and Intervention Controls (Simmhan et al., 2005). Each autonomous decision is associated with a Record of State Conditions; Evaluated Actions; Considered Actions; and Enforced Constraints (van der Aalst, 2016). Such a record supports Auditability and facilitates Post-Execution Analysis of System Behavior (Raji et al., 2020). Oversight also involves mechanisms by which Humans can Intervene When Decisions Exceed Predefined Risk Thresholds, or When Uncertainty Exists Regarding Autonomous Execution (Mitchell et al., 2019). Oversight is designed to Operate at Appropriate Temporal Scales to allow Human Judgment to Supervise Patterns and Policies Rather Than Individual Actions (Papagiannidis et al., 2024). Designing Oversight to Operate at Strategic Levels Rather than Operational Levels allows Human Judgment to Align with Strategic Governance Rather than Micromanage Operational Activities (Laato et al., 2022).

The Integration of Governance with Execution Differentiates This Reference Architecture from Conventional Control Tower Or Automation Frameworks (Raji et al., 2020). Traditional Systems Often Rely on External Audits, or Periodic Reviews to Ensure Compliance (Laato et al., 2022). The Proposed Architecture Embeds Governance Within the Decision Loop Itself (Papagiannidis et al., 2024). The Integration of Governance and Execution Transforms Trustworthiness from an Aspirational Attribute to Measurable System Property (Mitchell et al., 2019). Autonomy, Trust, and Governance Are Not Competing Objectives but Mutually Reinforcing Elements Realized Through Architectural Design (Raji et al., 2020).

This Reusable Reference Architecture Provides a Blueprint That Can Be Adapted Across Industries and Operational Contexts (Koot et al., 2021). The Layered Structure Allows for Modular Implementation Allowing Organizations to Incrementally Adopt Agentic Capabilities While Maintaining Governance Continuity (Monostori et al., 2016). By Formalizing the Relationships Among Data Perception; Decision Intelligence; Execution; and Oversight the Architecture Advances Supply Chain Research Beyond Isolated Algorithmic Contributions (Abideen et al., 2021). The Architecture Establishes a Systems Level Foundation for Trustworthy Agentic Supply Chains Able to Operate Responsibly in Complex Regulated and High-Consequence Environments (Papagiannidis et al., 2024).

Figure 2: Reference Architecture for Trustworthy Agentic Supply Chains



A Layered Reference Architecture for Trustworthy Autonomy within Agentic Supply Chains integrates Data Perception, Decision Intelligence, Execution, and Governance in a Single Controlled Decision Loop. The layers are designed to be vertically integrated to provide a complete picture of each component of the architecture. The first layer, Data Ingestion and Contextualization, is at the base of the architecture. It captures real-time signals from enterprise systems (e.g., ERP, WMS, TMS), other data sources (e.g., telemetry, partner communications, logistics events, regulatory data) and normalizes, temporally aligns, estimates confidence and detects anomalies to create an uncertainty aware system state. The Data Ingestion and Contextualization layer sends its output to the second layer, Decision Intelligence, which has four main components: Forecasting, Optimization, Learning, and Policy Generation. The Decision Intelligence layer continuously generates executable agent policies based on organizational objectives and constraints and not simply static recommendations. These policies are then sent to the third layer, Agent Execution and Coordination. Agents in the Agent Execution and Coordination layer are domain specific (inventory allocation, transport orchestration, etc.) and have limited autonomous capabilities; they coordinate with each other using their common policy context to generate validated commands to physical systems. All valid commands produced by the agents are sent through the fourth layer, Action Validation and Contextualization, to ensure compliance with organization objectives, constraints, and risk thresholds before those commands are executed. The Action Validation and Contextualization layer produces audit records and logs of all valid actions taken by the agents and safeguards these logs to support future traceability and accountability. Finally, the fifth layer, Guardian Agent and Governance, is the overarching layer of the architecture. It receives trust signals from system behavior, audit evidence, and performance outcomes and supports human oversight by providing dashboards, reports, escalation paths and intervention controls. While control flows are primarily top down in the form of policies, constraints, and authority boundaries, feedback flows are primarily bottom up in the form of telemetry, outcomes, and audit signals and as such creates a closed-loop system in which autonomy, trust, and governance are not external add-ons but part of the architecture itself.

Agentic Decision Intelligence and Levels of Autonomy

Agentic decision intelligence in the area of supply chains requires a specific determination of how autonomy will be assigned to decision-making areas (scopes) and organization levels (Parasuraman et al., 2000). Without such assignment, autonomous systems run the risk of either not performing adequately because they were overly constrained; or autonomously overstepping to the point that they violate the principles of stability, compliance, and accountability. A structured method of determining the degree of artificial decision-making capabilities in relation to organizational intent has been identified as "levels of autonomy" (Parasuraman et al., 2000). Levels of autonomy recognize that the amount of autonomy that should be given to artificial agents varies significantly among decisions regarding the degree of independence required, temporal urgency, and level of risk tolerance. The formalization of these distinctions is necessary for establishing controlled autonomous decision-making, as opposed to permitting uncontrolled autonomy and enabling meaningful operational adaptations.

There is a basic distinction between task-level autonomy and system-level autonomy (Scerri et al., 2002). Task-level autonomy applies to decisions related to a narrow scope of tasks that operate within defined boundaries (e.g., adjusting reorder quantities; selecting carriers; redistributing inventory across proximate nodes). These types of decisions are usually repetitive, time-sensitive, and have clearly defined constraints. Artificial agents are able to execute these decisions rapidly and consistently without requiring continuous human intervention. On the other hand, system-level autonomy involves decisions that affect the structure of the supply chain as a whole (e.g., configuring the supply chain; determining sourcing strategies; investing in capacity). These decisions are characterized by a higher degree of uncertainty, a longer time horizon, and higher organizational risk. Therefore, agentic architectures need to determine which decisions are capable of being executed autonomously by artificial agents, and which decisions require human authorization or collaborative oversight.

In addition to providing a distinction between task-level and system-level autonomy, the distinction between these two types of decisions also represents a difference in terms of decision coupling and consequences of decision-making (Klein et al., 2004). Decisions made at the task level generally result in localized effects that can be corrected or reversed with minimal impact to the rest of the supply chain. Conversely, decisions made at the system level may change the possible courses of action available to multiple subsequent processes, and create path dependencies that are difficult to reverse. Organizations can utilize this understanding to develop autonomy

regimes that allow artificial agents to make decisions where those decisions offer the greatest value-added, while maintaining human oversight of decisions that have strategic or ethical implications. This layered approach to agentic intelligence and organizational governance avoids replacing one with the other.

Another fundamental principle in agentic decision intelligence is the distinction between policy learning and policy execution (Schulman et al., 2017). Policy learning relates to the process by which an artificial agent learns to update its internal decision logic, based upon the results of past decisions; the receipt of new data; or changes in the supply chain environment. Policy execution, conversely, relates to the actual execution of decisions that ultimately affect the physical supply chain. Combining these functions creates instability, since exploratory learning behaviors may negatively affect operations. Separating learning from execution provides assurance that adaptive improvements to decision policies are both evaluated and validated prior to being deployed into live operational environments. This separation also allows learning to continue continuously, while ensuring that execution remains predictable and governable.

Separating learning from execution also facilitates organizational accountability (Mitchell et al., 2019). Policies developed through the learning process are subject to review, testing, and validation within existing organizational governance structures prior to expansion of execution authority. This process mimics the standard practice in safety-critical systems, in which control logic is verified and validated prior to deployment. In agentic supply chains, the separation enables the use of digital twins for evaluating the performance of learned policies in simulated scenarios (Mitchell et al., 2019). Execution authority is only granted once learned policies have demonstrated satisfactory performance, stability, and compliance. This design principle converts learning into a controlled mechanism for ongoing improvement, instead of a potential source of risk.

Bounded autonomy is the mechanism for establishing operational limitations on the decision authority provided to artificial agents (Altman, 1999). Instead of providing artificial agents with complete freedom of action, bounded autonomy defines specific limits on the scope, magnitude, and type of actions that can be taken by agents. These limits may take the form of quantitative thresholds (e.g., the maximum amount of inventory that can be reallocated), temporal constraints (e.g., the maximum number of decisions that can be made per unit of time), or qualitative constraints (e.g., the suppliers from whom orders cannot be placed; the geographic regions in which the agents are not allowed to operate). The primary purpose of bounded autonomy is to ensure that artificial agents continue to act in accordance with the organizational risk tolerance and regulatory requirements, regardless of the degree of autonomy that is granted to them. Importantly, the bounds applied to artificial agents' decision-making authority do not remain static, but rather are dynamically modified as the confidence in the ability of the agents to make decisions increases, or as the operating conditions of the supply chain change.

In addition to recognizing the importance of establishing bounds on artificial agents' decision-making authority, the concept of bounded autonomy also acknowledges that decision-making in supply chains occurs under uncertainty and with incomplete knowledge (Gu et al., 2024). By limiting the magnitude of actions that can be taken by artificial agents when the degree of uncertainty is high, the architecture ensures that artificial agents do not make aggressive decisions that exacerbate volatility. The bounds function as stabilizers that adaptively modulate the degree of autonomy provided to artificial agents in response to changing levels of confidence in their decision-making abilities and/or changes in the supply chain environment. As a result, agentic systems are able to maintain their effectiveness during normal operating conditions, while minimizing the degree of autonomy provided to artificial agents during periods of disruption or ambiguity. Bounded autonomy therefore provides a balance between the flexibility of artificial agents to respond to changing circumstances and the prudence of limiting the degree of autonomy provided to artificial agents during periods of uncertainty.

Finally, escalation thresholds provide a formal mechanism for transitioning decision authority from artificial agents to human overseers when predetermined conditions are met (Kaber & Endsley, 1997). These conditions may relate to the potential financial impact of a decision; the likelihood of triggering regulatory exposure; the risk of compromising safety; or the degree of deviation from expected behavior exhibited by the artificial agent. When a condition is reached, the artificial agent must temporarily suspend autonomous decision-making and seek human approval prior to resuming autonomous decision-making. This transition mechanism ensures that decisions that are extraordinary or have significant consequences are subjected to appropriate levels of scrutiny,

while continuing to permit artificial agents to make routine decisions without delay. Therefore, escalation thresholds allow human judgment to be preserved in situations where it is most relevant, without creating operational bottlenecks in day-to-day decision-making.

Escalation can formally be modeled as conditional control rule controlling agent actions (Ross & Varadarajan, 1989). Let a be a suggested action and r be the risk measurement related to action a , derived from system state and future outcome predictions. Escalation can formally be described as follows:

$$\text{Execute}(a) = \begin{cases} 1 & \text{if } r \leq \tau \\ 0 & \text{if } r > \tau \end{cases}$$

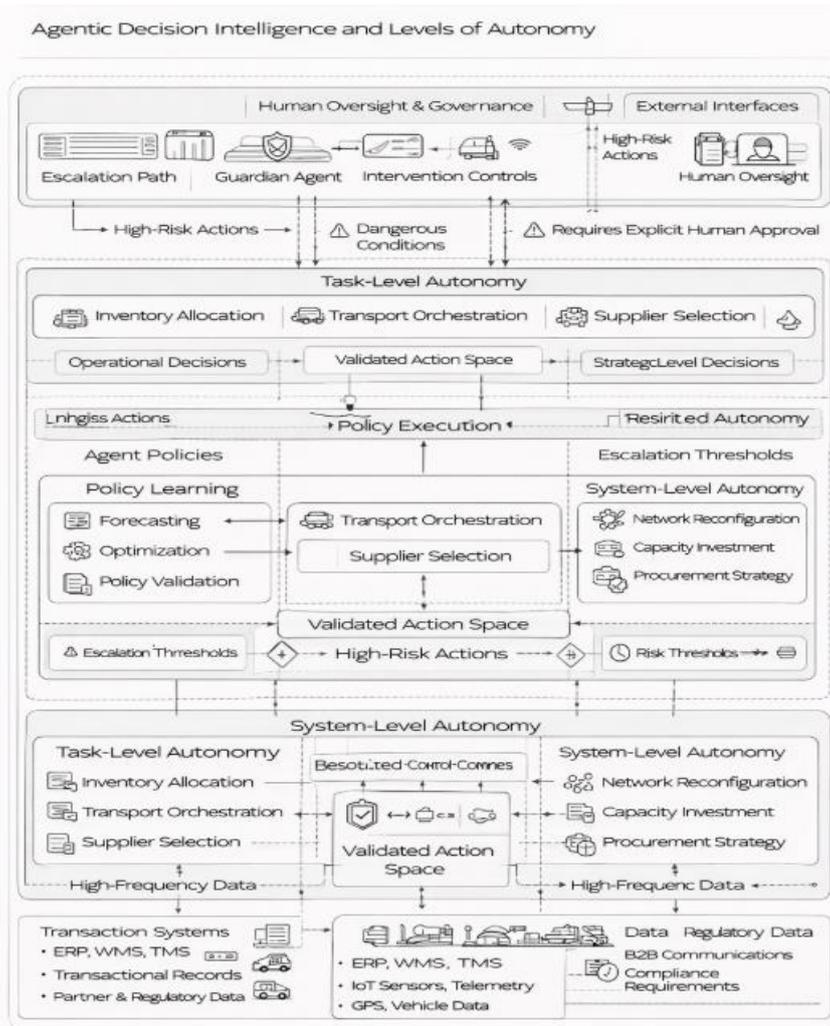
Where τ represents the escalation limit. This formulation clearly shows that the execution authority is dependent upon the evaluation of the risk, instead of simply depending upon the expected values. The formulation describes the architectural principle that autonomy is conditional and reversible based upon governance principles. Although the above formulation simplifies the representation of the logic through which agentic systems remain controllable and accountable, it does capture the logic that enables systems comprising autonomous agents to be collectively governed by common constraints and objectives.

In addition to enabling the decision-making capabilities of each agent, agentic decision intelligence also requires coordination among multiple agents that operate at various levels of autonomy (Amato, 2024). Task-level agents operate independently within their respective localized domains of operation, while system-level oversight agents or human supervisors ensure compliance with higher-order objectives. Mechanisms for coordinating agents include sharing a common state representation, establishing hierarchical control relationships, and establishing policy harmonization rules. These mechanisms prevent the occurrence of conflicting actions and ensure that local optimizations do not detract from overall system performance. Therefore, coordination enables autonomy to transform isolated decision-making into a collective capability governed by shared constraints and objectives.

The formalization of autonomy levels has important implications regarding trustworthiness (Papagiannidis et al., 2024). Trust in agentic systems does not stem exclusively from performance metrics; trust arises from confidence that decision authority is exercised appropriately and predictably. The architecture establishes the boundaries of autonomous authority by defining who can make decisions about what under which conditions. This legibility supports both internal organizational trust and external regulatory confidence. Stakeholders can understand the scope of authority that is granted to individual agents, and the safeguards that restrict the exercise of that authority, regardless of how quickly and autonomously decisions are made by machines.

Therefore, by structuring agentic decision intelligence using separate task and system level autonomy, separating learning and execution, bounding autonomy and establishing escalation limits, this framework prevents uncontrolled autonomy while allowing for continuous adaptive decision-making to continue (Li & Goel, 2025). Furthermore, autonomy is not viewed as an end-state, but rather as a calibrated capability that is situated within governance structures. This approach views agentic supply chains as controllable adaptive systems whose intelligence improves organizational decision-making, without undermining accountability or institutional legitimacy.

Figure 3: Agentic Decision Intelligence and levels of Autonomy



The diagram illustrates a technically layered architecture for agentic decision intelligence that embodies multiple levels of autonomy while maintaining governance, safety, and human control. The architecture consists of three layers: the lowest layer includes enterprise transaction systems, operational data sources, and other data feeds, including ERP, WMS, TMS, IoT sensors, telemetry, and regulatory data that provide high-fidelity inputs to the architecture that define the real-time state of the supply chain. The next layer includes the decision intelligence layer that utilizes these input data to generate agent policies through the use of forecasting, optimization, learning, and policy generation modules. Importantly, this layer is responsible for generating policy through learning alone and not through direct execution, and thus allows for adaptive updates to be generated in a controlled manner. Once learned policies have been generated, they are forwarded to the uppermost layer, the agent execution and coordination layer, where autonomy is exercised within a validated action space. In this layer, task-level agents are responsible for making operational decisions, including, but limited to, inventory allocation, transport orchestration, and supplier selection. However, prior to exercising decision authority, proposed actions are evaluated against higher order objectives, constraints, and quantifiable risk and cost measures. The validation process exercises bounded autonomy through limiting the scope and magnitude of permissible actions and their frequencies and prevents agents from taking decisions that contravene organizational or regulatory limits. Additionally, this layer embeds escalation thresholds that continually assess decision risk and route control to the uppermost layer through explicit escalation pathways whenever decisions are identified as being at high risk or anomalous. The uppermost layer of the architecture is comprised of the human oversight and governance layer that includes guardian agents, dashboards, intervention controls, and external interfaces. This layer maintains decision authority over system-wide autonomy and is responsible for making strategic decisions, such as network reconfigurations, capacity investments, and procurement strategies, that require explicit human approval and are thereby excluded from autonomous execution. Through-out the architecture, the separation between policy learning and execution, along with the incorporation of audit logs,

safeguarded records, and traceable escalation mechanisms, ensures that autonomy remains conditional, reversible, and accountable. From a technological standpoint, the diagram encodes autonomy as a constrained control hierarchy in which agents optimize within delegated bounds, governance rules establish possible action sets, and human oversight intervenes when the risk of an action exceeds pre-defined thresholds, thereby integrating trust, adaptability, and institutional control into a single executable system design.

Governance by Design in Agentic AI Systems

Embedded Governance Constraints

Embedded governance constraints provide the structural basis for autonomous decision-making systems in supply chain systems to address issues related to organizational authority and institutional accountability as they execute autonomous decision-making. The increasing use of artificial agents in global supply network operations has resulted in many of the decisions associated with sourcing allocation, routing, inventory positioning, and fulfillment occurring at timescales that exceed the capabilities of human decision makers (Parasuraman et al., 2000). Therefore, as artificial agents make operational decisions, governance cannot be solely an activity that occurs post-execution (Papagiannidis et al., 2025) but instead must be incorporated internally into the system architecture itself. Constraints ensure that all autonomous decisions are made only within the bounds of what the organization is permitted to do and what the organization is willing to do; therefore, governance is transformed from a corrective process into a component of decision intelligence (Floridi & Cowsls, 2019).

In supply chain systems, governance constraints illustrate a complex intersection of trade regulations, contractually agreed terms, risk management policies, ethical commitments, and strategic priorities. Trade regulations may prohibit sourcing from a particular region. Financial governance may place restrictions on inventory levels or work capital utilization (Altman, 1999). Sustainability commitments may establish limitations on greenhouse gas emissions or require suppliers to meet certain certifications (Ivanov & Dolgui, 2021). When constraints are embedded in the same decision environment as the agent making the decision, the artificial agent does not view them as separate, external check processes but as inherent components of the feasible action space (Altman, 1996). The incorporation of constraints into the decision-making process will ensure that each autonomous decision made will be aligned with the organizational obligation at the time the decision is made, not subsequent to the decision (Altman, 1996).

The absence of embedded governance provides structural risk to autonomous supply chain systems. Learning-based agents will seek to optimize their objectives based on reward functions (Sutton, 1988). If constraints are not embedded in the reward functions, agents may discover decision strategies that result in optimized objective functions (short-term) and violate regulatory or ethical constraints (Alshiekh et al., 2018). In supply chain systems, this may take the form of extreme cost reduction due to exploiting regulatory arbitrage, unsafe labor practices, or unsustainable routing options (Garvey et al., 2015). Although a post-hoc compliance review may identify the violation(s) after damage has occurred, it cannot prevent financial penalties, reputational harm, or disruption to supply. Embedded governance will eliminate this risk by preventing the artificial agent from exploring regions of the state-action space that are prohibited (Brunke et al., 2022).

From a systems theory perspective, embedded governance constraints serve as invariant boundary conditions on agent behavior (Chow et al., 2018), defining areas of the state-action space that will be categorically excluded regardless of changes in environmental uncertainty or performance pressure (Chow et al., 2018). These types of boundary conditions are particularly important in supply chains due to the propensity for disruptions (such as geopolitical conflict, natural disasters, etc.) to encourage extreme decision-making (Garvey et al., 2015). By embedding constraints such as trade compliance, safety, and other non-negotiable constraints into the decision-making process, the artificial agent's autonomy will not degrade into opportunistic decision-making during stressful periods (Lazarus et al., 2020). Maintaining the structural stability of autonomous decision-making is a necessary condition for the deployment of autonomous systems in mission-critical supply chains.

Embedded governance also transforms how artificial agents perceive and reason about the supply chain environment (Doshi Velez & Kim, 2017). Artificial agents do not have a neutral perception of the supply chain environment, but rather the environment is perceived through a filter that encodes the organization's priorities.

For example, inventory availability will be perceived differently when some inventory cannot legally cross borders or when contractual agreements reserve capacity for specific customers. By embedding governance into the way artificial agents construct states, artificial agents will reason about the supply chain environment based on a representation of the environment that is institutionally valid, not physically comprehensive (Doshi Velez & Kim, 2017). This distinction is important because artificial agents' decisions are only as trustworthy as the representations upon which those decisions were made (Doshi Velez & Kim, 2017).

The business benefits of embedded governance constraints extend beyond compliance into strategic risk management (Garvey et al., 2015). Supply chains operate under conditions of asymmetry of information and delayed feedback (Garvey et al., 2015). Decisions that appear optimal in one region of the state-space may create systemic risk in another region of the state-space over time due to the accumulation of exposures or cascading failures (Garvey et al., 2015). Embedded constraints allow organizations to proactively enforce risk limits by limiting decision patterns that increase fragility, even though they may appear beneficial individually (Garvey et al., 2015). Examples of constraints on decision patterns include supplier concentration and transportation lane dependencies. Organizations can algorithmically enforce such constraints to preserve long-term resilience in exchange for some potential short-term inefficiency. This allows organizations to link autonomous execution with enterprise-wide risk governance objectives.

Embedded governance constraints also support scalability in global supply chain operations. As organizations grow globally, the number of regulations and contracts that must be complied with grows exponentially (Altman, 1999). Human oversight alone cannot effectively monitor the vast numbers of autonomous decision-making executions that occur in supply chains on a daily basis. By incorporating jurisdiction-specific rules into the decision-making architecture of artificial agents, organizations can scale the governance enforcement of autonomous decision-making without scaling the managerial workload of monitoring autonomous decision-making. This scalability is an important enabler for multinational enterprises that wish to deploy agentic supply chains in multiple regulatory environments while maintaining consistent governance standards.

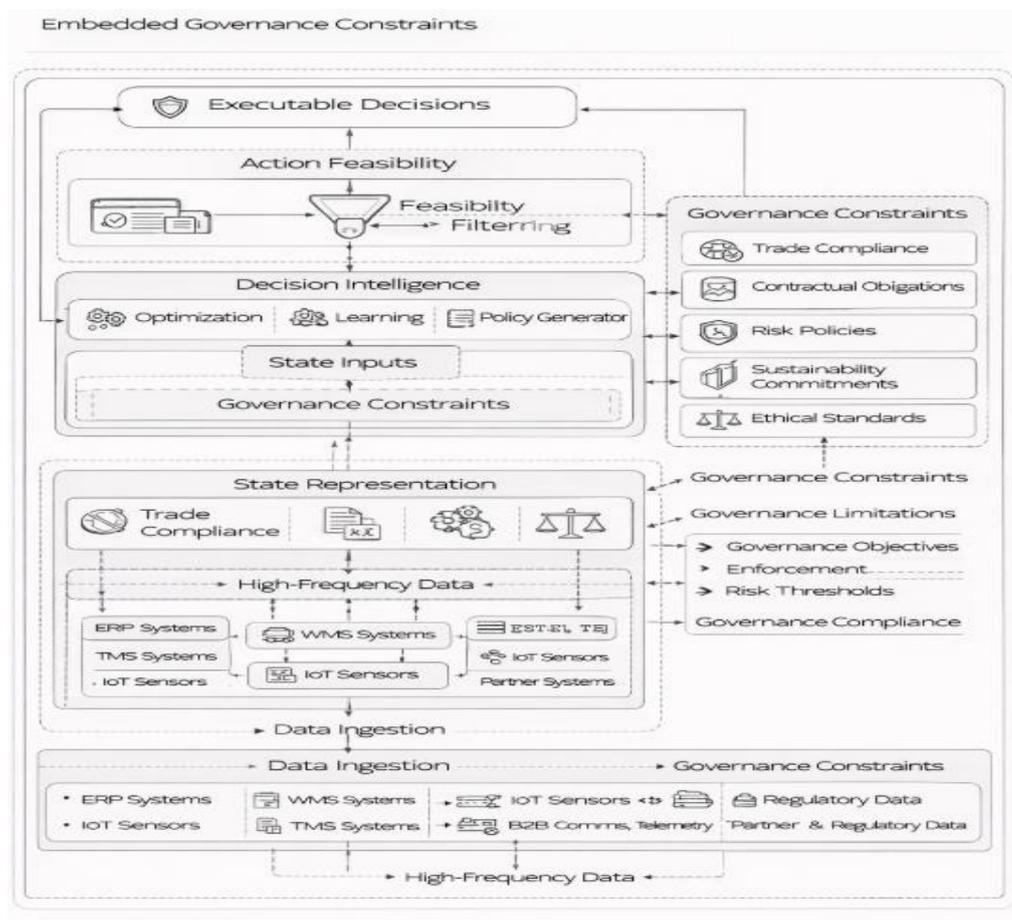
Another theoretical implication of embedded governance relates to the relationship between governance and learning (Brunke et al., 2022). In unconstrained learning systems, agents may explore actions that are institutionally unacceptable (Brunke et al., 2022). Embedded constraints limit exploration to acceptable regions of behavior and ensure that learning trajectories remain aligned with organizational norms (Alshiekh et al., 2018). Such limitation does not weaken intelligence but rather focuses intelligence on viable strategies in practice. In supply chains, this accelerates convergence towards deployable policies since agents do not waste learning capacity on infeasible actions. Governance therefore increases the efficiency of learning while maintaining institutional integrity.

Embedded governance also supports transparency and accountability in autonomous supply chain systems (Raji et al., 2020). Since constraints are explicitly represented, stakeholders can examine and validate the normative assumptions underlying agent behavior (Raji et al., 2020). Transparency is important for stakeholder trust (regulatory agencies, partners, etc.), internal governance reviews, and regulatory engagement (Raji et al., 2020). Rather than relying on informal assurances that systems behave responsibly, organizations can demonstrate that responsibility is built into the structure of the system (Raji et al., 2020). Such a demonstration of responsibility will strengthen organizational legitimacy and reduce barriers to the adoption of autonomous systems among stakeholders.

From an organizational design perspective, embedded governance constraints change the role of managers. Managers move from approving individual decisions to designing and implementing the governance structures that guide the behavior of autonomous decision-making systems (Klein et al., 2004). This represents an elevation in managerial responsibilities from operational interventions to architectural stewardship (Klein et al., 2004). In supply chains, this allows managers to focus on strategic alignment, policy coherence, and risk posture, while delegating autonomous decision-making to agentic systems. Embedded governance therefore provides the connection between human authority and machine autonomy.

In conclusion, embedded governance constraints transform agentic supply chains from technical autonomous systems into institutionally reliable operational infrastructures (Novelli et al., 2024). By embedding compliance, risk tolerance, and ethical commitments into the decision architecture of autonomous systems, governance is transformed from a reactive safeguard to a proactive determinant of behavior (Floridi & Cowls, 2019). Such a transformation is required to realize the business benefit of autonomy while minimizing the risk of regulatory and organizational risk (Novelli et al., 2024). Embedded governance provides the necessary infrastructure for autonomous systems to operate responsibly and maintain control legitimacy and strategic coherence.

Figure 4: Embedded governance constraints



Rule Constrained Policy Spaces

Embedded Governance Constraints are converted into formal Autonomy in Agentic Supply Chains via Policy Spaces. Embedded constraints specify what is institutionally acceptable in broad terms whereas Policy spaces constrain how embedded constraints limit the scope of the Actions Artificial Agents can execute. With many State Space possibilities and Combinatorial Complexity Decision-making in Agentic Supply Chain Environments, Unconstrained Policy Learning can lead to Artificial Agents executing Actions that are Technically Feasible but Operationally Not Acceptable. Constrained Policy Spaces resolve this Tension by limiting the Domain of Artificial Agent's Learning and Execution to Actions that are both Operationally Effective and Institutionally Legitimate.

Policy Spaces in Agentic Supply Chain Decision Making include Actions such as selecting Suppliers; allocating inventory; routing; prioritizing; scheduling Production; and Sequencing Fulfillment. Each Action exists in a Dense Web of Contractual Obligations; Regulatory Requirements; Capacity Limits; and Strategic Commitments. Thus, Artificial Agents learn how to Optimize their Performance Metrics within the Realities of Business Governance when Constrained Policy Spaces encode Boundaries of Governance Rules into the Mapping between States and Actions. This distinction is Critical to converting Theoretical Autonomy into Deployable Enterprise Systems.

Constraining Policy Spaces theoretically reconciles Adaptive Learning with Organizational Control by Reframing the Optimization Problem. Instead of Optimizing Value across All Possible Actions, Artificial Agents Optimize Value across a Restricted Set of Actions that reflects Governance Rules (Altman, 1999). Thus, Intelligence Emerges Only within Acceptable Boundaries. In Agentive Supply Chains, this is particularly important since Optimization Pressures are often at Odds with Compliance Obligations. For example, Cost Minimizing Strategies May Favor Suppliers with Regulatory Risk or Routing Paths with Geopolitical Exposure. Constraining Policy Spaces Ensures that Such Options are Excluded prior to Optimization Begins (Altman, 1996).

The Business Impact of this Approach includes its Capability to Prevent Governance Violations while Sustaining Operational Agility. Traditional Compliance Mechanisms typically Rely on Manual Approvals or After-the-Fact Audits, both of which Introduce Latency and Costs. Constraining Policy Spaces Algorithmically Converts Compliance into Automatic and Continuous Processes. Autonomous Decisions can be Executed at Scale without incurring Additional Oversight Burden. This Capability enables Enterprises to deploy Agentive Supply Chains Confidently Across Regions with Heterogeneous Regulatory Environments while Maintaining Consistent Governance Standards.

Constraining Policy Spaces also Influences Learning Efficiency and System Stability. By Eliminating Actions that are Infeasible or Prohibited from Consideration, the Effective Search Space for Learning Algorithms is Reduced. This Reduction Accelerates Convergence toward Effective Policies and Decreases the Likelihood of Erratic Behavior during Learning Phases. In Agentive Supply Chain Environments where Decision Instability can Propagate Rapidly through Interconnected Networks, this Stabilizing Effect has Direct Financial and Operational Benefits. Smoother Convergence Reduces Experimentation Costs while Preserving Service Reliability.

Dynamic Supply Chain Environments may have Varying Governance Rules depending upon Context. Constraining Policy Spaces Support Conditional Constraints that Activate Based upon State Conditions. For Example, Routing Options may be Permissible under Normal Conditions but Restricted during Periods of Geopolitical Tension. Inventory Transfers may be Allowed within Certain Regions but Prohibited across Jurisdictions due to Data or Trade Restrictions. By Incorporating such Conditional Logic into Policy Spaces, Agentive Systems can Adapt their Behavior Dynamically while remaining Compliant. This Contextual Flexibility Enhances Resilience without Undermining Governance Integrity.

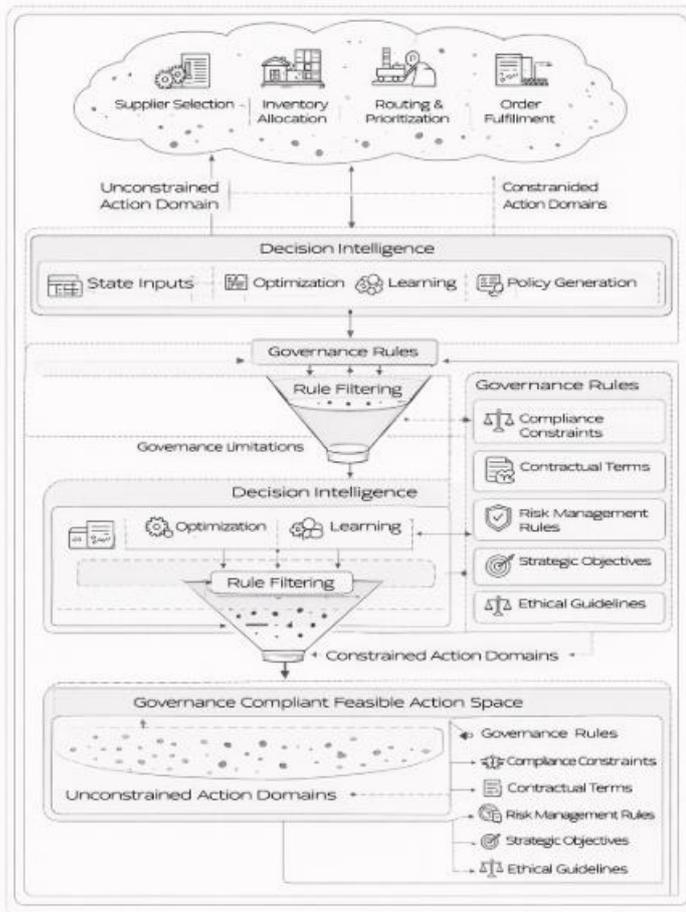
Constraining Policy Spaces also Support Transparency and Auditability in Autonomous Decision Systems. When Actions are Selected from Explicitly Defined Permissible Sets, Decision Rationales become Easier to Reconstruct and Evaluate. Auditors and Regulators can Assess not only the Outcome of a Decision but whether the Action was Allowable under Governing Rules at the Time of Execution. This Clarity Strengthens Institutional Trust and Simplifies Compliance Reporting. In Contrast, Unconstrained Systems Require Complex Post-Hoc Explanations to Justify why Certain Actions were Taken.

Organizational Governance Perspectives view Constraining Policy Spaces as a means to Transfer Responsibility from Individual Decision Approval to Rule Definition and Maintenance. Leaders and Compliance Teams Define the Boundaries of Acceptable Behavior while Agents Operate Independently within those Boundaries. This Separation enables Governance Expertise to be Concentrated Where it Adds the Most Value while Execution Scales Algorithmically. The Result is a Governance Model that is both Rigorous and Efficient, aligning with the Scale and Complexity of Modern Supply Chains.

Theoretical Rigor is Required to Ensure that Rule Constraints do not Conflict or Create Infeasible Decision Regions. Poorly Designed Constraints may unintentionally eliminate all viable Actions under certain States, leading to Escalation or System Paralysis. Addressing this Risk Requires Systematic Validation of Policy Spaces through Simulation and Stress Testing. In Agentive Supply Chains Digital Twins play a significant Role in Evaluating how Constrained Policies Behave under Extreme Conditions, ensuring that Governance Rules Preserve Feasibility while Enforcing Compliance (Fuller et al., 2020).

Therefore, Constraining Policy Spaces represents a Central Pillar of Governance by Design. They Transform Governance from a Monitoring Function into a Generative Force that Shapes Intelligent Behavior. By Defining the Space within which Autonomy Operates, they Enable Agentic Supply Chains to Achieve Adaptive Performance without Compromising Institutional Obligations. This Capability is Essential for Realizing Business Value from Autonomy while Preserving Regulatory Compliance and Organizational Legitimacy.

Figure 5: Rule Constrained Policy spaces



Human on the Loop Supervision

Human-on-the-loop (HOTL) supervision facilitates the alignment of strategic intent and tactical operation in autonomous systems that operate at machine speed (Parasuraman et al., 2000). Autonomous systems cannot require continuous approval for every decision made individually as it would not be possible given the sheer volume of decisions made by supply chain systems hourly by day for thousands of hours. However, HOTL supervision addresses the scale issue by transforming how human oversight occurs from an operational role to a systemic role, and thus, humans do not need to approve every single action, rather they review the policies and outcomes of the autonomous decisions being made by the systems (Kaber & Endsley, 1997).

The theoretical base for HOTL supervision is the differentiation of decision governance and decision execution. Decision execution involves identifying and implementing responses to immediate situations. Decision governance includes ensuring that the processes used to make decisions remain consistent with the organization's objectives, risk profile, and ethics. Autonomous systems have the ability to execute decisions quickly and consistently; however, humans have the ability to reason ethically, interpret contextually, and assume accountability for those decisions. Therefore, HOTL supervision formalizes this separation of duties between humans and autonomous systems within supply chain architectures.

In actual supply chain operations, human supervisors will oversee aggregate measures including service level trend analysis, cost variances, risk exposures, compliance and policy drifts. Those aggregate measures allow for

the determination if the autonomous behavior is remaining within the desired boundaries. Rather than intervene at each decision point, supervisors will modify the parameters of governance, add constraints or alter escalation points when trends emerge indicating the autonomous behavior is not remaining in line with organizational goals (Bryson & Crosby, 1992). The business impacts of HOTL supervision are numerous. Organizations that remove humans from the approval process for routine decisions improve the responsiveness and efficiency of their supply chain operations. Additionally, by allowing for human oversight to occur at the policy level, organizations eliminate the reputational and regulatory risks associated with having autonomous systems without some form of oversight (Morgeson et al., 2015).

Additionally, HOTL supervision supports adaptive governance in dynamic environments. Supply chains are subject to sudden and unpredictable changes in markets, geopolitics, and regulatory environments. Human supervisors can respond to changes in these areas by adjusting governance rules or modifying the boundaries of autonomy in almost real-time. These adjustments then cascade throughout the system architecture and reshape the autonomous behaviors of the agents in the system, all without shutting down the system or manually intervening (Camarinha-Matos et al., 2016). This capability improves the resiliency of supply chains by providing governance that evolves in tandem with the changing environment.

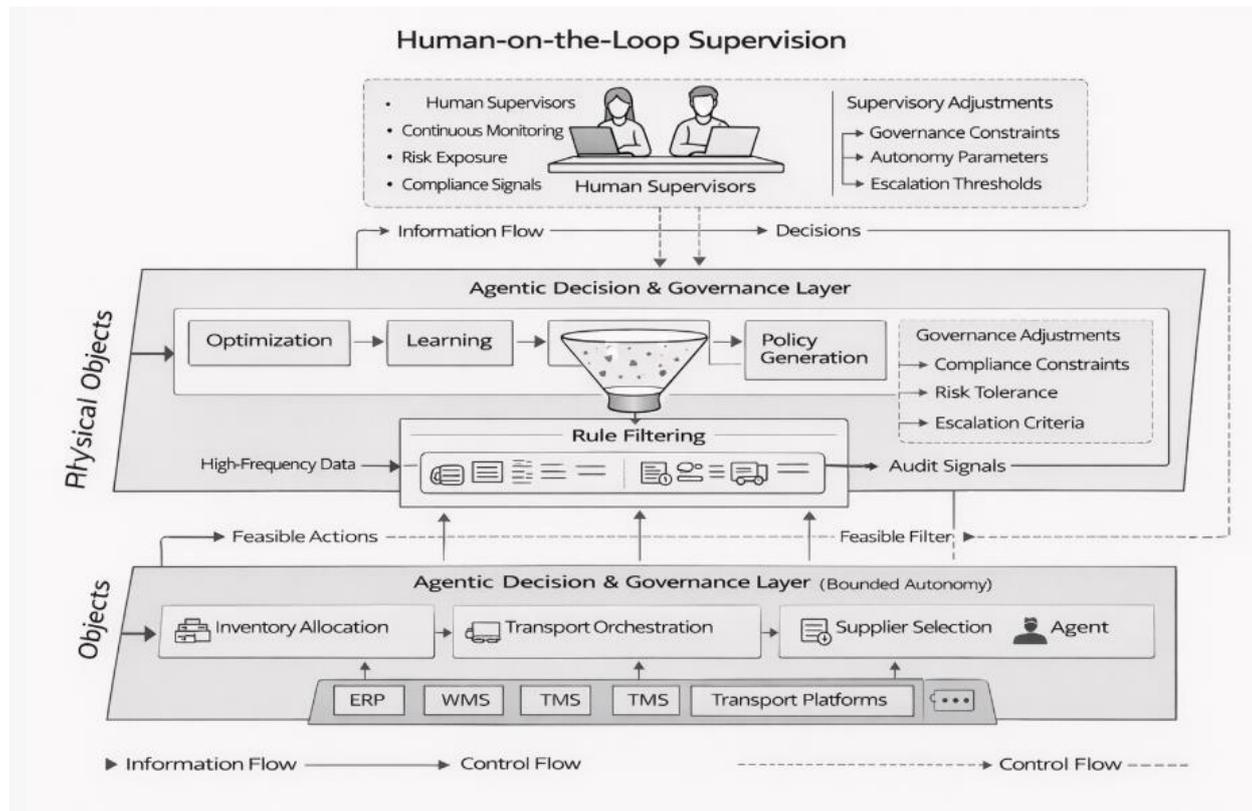
From a cognitive standpoint, HOTL supervision recognizes the limitations of human attention and decision-making. Humans are poorly suited to monitor the frequency of high-frequency events, but are exceptional at recognizing patterns and anomalies across aggregated data. Thus, by presenting supervisors with synthesized indicators, rather than raw decision streams, the system enables effective oversight without placing cognitive demands upon the supervisor. This design leverages the strengths and weaknesses of both humans and machines. Lastly, HOTL supervision provides for accountability in autonomous supply chains. Agents carry out decisions, while humans retain responsibility for developing the governing framework for the decision-making process. This is a fundamental principle for establishing accountability for both legal and ethical purposes. Organizations can establish that autonomous behavior has not been left ungoverned, but has occurred under continuous human oversight at the policy level. Establishing this accountability structure is essential for gaining regulatory acceptance of autonomous systems in high-risk environments.

To ensure that the interaction between human supervisors and autonomous systems is truly bidirectional, supervisors receive information regarding the behavior of the system, and in turn, supervisors have the opportunity to influence future behavior via adjustments to governance. This feedback loop ensures that autonomy remains responsive to the organizational learning and external expectations. As supervisors become familiar with the performance and transparency of the system over time, they build trust and confidence in the system, allowing for a gradual increase in the degree of autonomy granted to the system, as necessary.

Finally, HOTL supervision is a key component of management of organizational change. The implementation of autonomous systems affects the roles, responsibilities, and decision authority of individuals throughout the organization. The supervisory interfaces enable managers to maintain involvement with supply chain operations without becoming mired in operational detail. Maintaining this involvement supports the adoption of autonomous systems, and maintains the perception of control and understanding throughout the transition to autonomous operations.

Overall, HOTL supervision transforms autonomous systems from a threat to managerial authority to an extension of organizational capabilities. By transitioning human oversight from the execution of decisions to the governance of decisions, autonomous supply chains can achieve scalability, adaptability, and accountability concurrently. This supervisory model is foundational to linking autonomous decision-making systems to business strategies, regulatory expectations, and ethical considerations.

Figure 6: Human-on-the-loop Supervision



Emergency Intervention Mechanisms

Emergency intervention mechanisms provide an ultimate safeguard for governance by design in agentic supply chain systems; the other safeguards – embedded constraints, rule-constrained policy spaces, and human-in-the-loop (human-on-the-loop) supervision – generally prevent governance failures, but complex adaptive systems will still occasionally face an event or scenario which is outside of what has been modeled. Supply chains have many examples of extreme events including: a global war; a sudden regulatory change; a cyber-attack; a major disaster (e.g., earthquake, hurricane); a cascade failure of suppliers; etc. These types of events are outside of the modeled range of the system, and emergency intervention mechanisms provide a means to reverse the autonomous decision-making process to preserve organizational control and institutional legitimacy when normal governance pathways fail.

Emergency intervention mechanisms, therefore, establish the boundary conditions of autonomy. The autonomy of agentic supply chain systems is not absolute, and is contingent upon the presence of environmental stability and acceptable levels of risk. The concept of emergency intervention formalizes the idea that decision authority can be withdrawn or reorganized when systemic risk exceeds pre-defined tolerance levels. The reversibility provided by emergency intervention mechanisms is fundamental to the establishment of trustworthiness since they guarantee that no autonomous system exists outside of human reclaim.

For those supply chains which make decisions affecting physical goods, financial exposure, and regulatory compliance, the ability to intervene quickly and decisively in extreme situations is a necessary condition to achieve widespread adoption.

Emergency intervention mechanisms function at multiple levels of the supply chain architecture. At the agent level, intervention can include suspending an individual agent(s) who exhibit behaviors significantly different than what were anticipated. At the system level, intervention can include halting specific types of decisions (i.e. cross-border transactions, supplier switching). At the organizational level, intervention can include the transfer of all decision authority to humans for a specified time frame. A multi-level design provides the ability to target

interventions in a manner that is proportional to the severity of the problem, and avoid blanket shutdowns that can create new operational risks.

Emergency intervention is inherently linked to Business Continuity Management (BCM) in supply chain operations. Autonomous systems are frequently used to improve responsiveness and resilience, however, uncontrolled autonomy in crisis situations can increase instability. For example, if there is a sudden collapse of demand, or a shock to supply, autonomous agents may react aggressively to rebalance inventory, cancel orders, etc. in ways that damage relationships with suppliers, or violate contractually obligated commitments. Emergency intervention mechanisms provide an organization the ability to stabilize behaviors through reverting to conservative base-line policies that emphasize continuity over optimization. This enables protection of long-term business relationships and brand reputation in times of stress.

An important design consideration of emergency intervention mechanisms is to promote graceful degradation, rather than abrupt termination. Termination of autonomous execution can cripple supply chain operations and create new risks such as order backlogs or missed shipments. Gradual degradation involves the controlled reduction of autonomy through the narrowing of the action boundaries, the raising of escalation thresholds, or the transition to slower, but safer control modes. This gradual response to emergency conditions allows the supply chain to continue to function while limiting the scope of risk. Additionally, from a business perspective, this approach reduces the degree of disruption, and promotes managerial confidence.

Emergency intervention mechanisms also mitigate the risk of model drift and systemic misalignment. Even well-governed agentic systems can suffer from drift due to changes in demand patterns, supplier behavior, and regulatory environments. When sufficient drift occurs, emergency intervention mechanisms allow organizations to stop execution and begin to recalculate. In supply chains, this prevents the reinforcement of maladaptive decision-making processes that could result in sustained inefficiencies or exposure to compliance issues. Thus, emergency intervention mechanisms serve as a corrective reset mechanism to ensure long-term viability of the system.

The determination of whether emergency intervention is required, uses composite risk indicators, rather than relying solely on one metric. Risk in supply chains emerges from the interaction of several factors including: inventory exposure; service volatility; financial commitments; and regulatory sensitivity. Emergency triggers are activated by the crossing of thresholds on composite risk measures rather than a single decision outcome. Therefore, this holistic approach to risk assessment is consistent with Enterprise Risk Management best practices, and ensures that emergency interventions are justified based on the systemic conditions, and not by transient anomalies.

An easy to understand, formalized version of emergency intervention logic is a System Level Risk Evaluation. We will define R_t = An aggregate risk metric at time t based on project operational and compliance metrics. Then we can represent intervention as:

$$I_t = \begin{cases} 1 & \text{if } R_t > \kappa \\ 0 & \text{if } R_t \leq \kappa \end{cases}$$

Where κ represents the Emergency Intervention Threshold (Altman, 1999). Once intervention is initiated, then either Autonomous Execution is suspended or limited based on predefined Protocols. This representation clearly defines Emergency Control as a response to a system-wide risk condition as opposed to a localized Optimization Failure. This definition further emphasizes Governance's priority of Organizational Survival/Legitimacy over Short Term Performance.

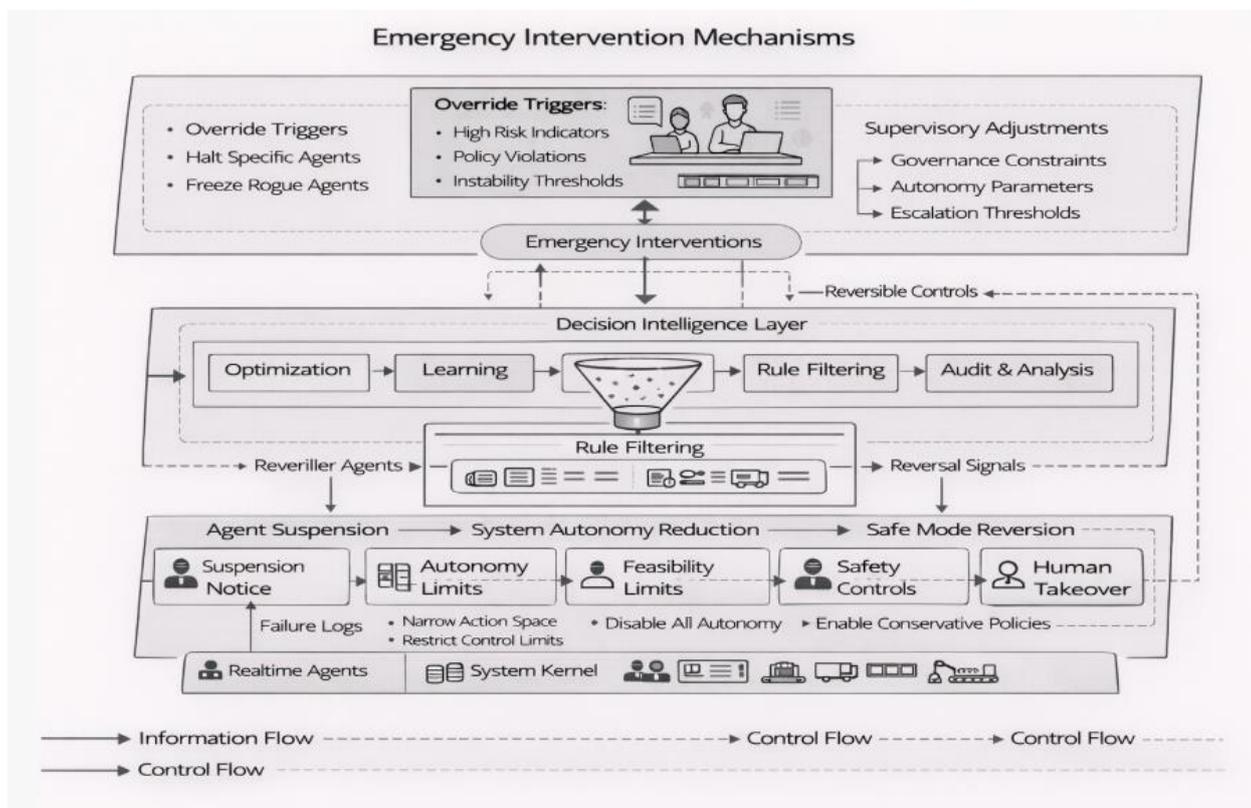
Regulatory Acceptance of Agentic Supply Chains (Schiff et al., 2024), rely heavily on Assurance that Autonomous Systems can be stopped or overridden when Compliance Risk arises. Tested Intervention Pathways provide assurance that Autonomy has been utilized Responsibly and remains subordinate to Legal Authority. Regulatory Approval in areas like Pharmaceutical, Food Distribution, Defense Logistics etc., requires Assurance that Autonomous Systems can be Stopped/Overridden. Therefore, Emergency Mechanisms allow for both Operational Safety and Regulatory Viability.

From an Organizational Governance Perspective Emergency Intervention Provides Clear Accountability for Executive Leaders. Executive Leaders maintain Ultimate Authority over Autonomous Systems and are authorized to intervene when Strategic/Ethical issues arise. As such, the Clarity of Authority reduces concerns that Autonomous Systems reduce Managerial Control. Instead, Autonomy is viewed as a Delegated Capability which can be Withdrawn Under Defined Conditions. This Viewpoint Supports Confidence in Executives and Facilitates Adoption within Conservative Organizational Cultures.

Post-Intervention Analysis generates Learning Opportunities that Enhance Governance Over Time. The Insights generated from post-intervention Analysis provide Understanding of Failure Modes Governance Gaps and Environmental Dynamics that were Not Fully Anticipated. Such Insights Inform Updates to Embedded Constraints Policy Spaces Escalation Thresholds and Supervisory Metrics. As such, Emergency Intervention Contributes to Organizational Learning as well as Representing a Failure of Autonomy. Supply Chains Evolve not Only Through Optimization but Through Disciplined Responses to Extreme Events (Ivanov & Dolgui, 2021).

In Summary, Emergency Intervention Mechanisms Complete the Governance by Design Framework by Providing Assurance that Agentic Supply Chain Systems Remain Controllable Under Extreme Uncertainty. By Allowing Targeted Suspension Graceful Degradation and Reversion of Authority these Mechanisms Protect Institutional Legitimacy While Preserving Operational Continuity. Emergency Intervention Transforms Autonomy From an Irreversible Commitment Into a Managed Capability That Can Be Exercised Confidently in High Consequence Supply Chain Environments.

Figure 7: Emergency Intervention Mechanisms



Compliance and Regulatory Alignment

Mapping Agentic Decisions to Regulatory Requirements

Autonomous Supply Chain Systems – Mapping Decisions to Requirements
Mapping agentic decision-making to regulatory requirements is a key factor in successfully deploying autonomous supply chain systems in both highly regulated and high-consequence settings (Raji et al., 2020). Prior to the advent of agentic supply chains, traditional compliance methodologies assumed that decisions were

made through distinct human actions that could be assessed post-implementation (Vasarhelyi et al., 2004). Agentive supply chains operate in an environment of tens-of-thousands of operational decisions being executed continuously, across sourcing, logistics, inventory management, and fulfillment (Alles et al., 2008). Therefore, regulatory risk in agentive supply chain systems is not typically due to the occurrence of singular decisions; instead it is the cumulative effects of repeated behaviors occurring over time (Böhmecke Schwafert, 2024). To ensure that agentive decision-making is aligned with regulatory requirements, it is necessary to map those requirements directly onto decision logic so that the autonomous nature of decision-making can remain aligned with the institutions' goals as those decisions occur rather than solely through retrospective enforcement mechanisms (Raji et al., 2020).

Regulatory Requirements of Supply Chains: Trade laws, environmental regulations, labor standards, financial reporting requirements, and data protection requirements provide examples of regulatory requirements for supply chains, however these are defined at an institutional level (e.g., international law) and do not provide clear guidance on how to evaluate individual routing, or sourcing decisions (Tong et al., 2022). As such, mapping agentive decisions to regulatory requirements involves translating normative regulatory rules into operational constraints that autonomous systems can interpret (Raji et al., 2020). This translation will bridge the gap between institutional intent and computational logic to enable regulatory considerations to influence decision-making in real-time (Böhmecke Schwafert, 2024).

Supply Chain Decision-Making Impacts Multiple Domains

In agentive systems, supply chain decisions create regulatory exposures across multiple domains concurrently (Tong et al., 2022). For example, a sourcing decision may require compliance with trade sanctions, labor practices, and sustainability requirements (Tong et al., 2022). Similarly, a routing decision may require compliance with customs procedures, safety standards, and emissions limits (Chalendar et al., 2019). Inventory allocation decisions may also create regulatory obligations related to tax jurisdiction and data localization (Böhmecke Schwafert, 2024). To effectively map regulatory requirements to agentive decisions, it is necessary to associate each decision domain with the applicable regulatory frameworks, and have agentive systems evaluate compliance in context (Raji et al., 2020). The mapping of regulatory requirements to agentive decision-making prevents regulatory blind-spots that result from evaluating decisions in isolation (Böhmecke Schwafert, 2024).

Increased Complexity of Global Supply Chains

The complexity of mapping regulatory requirements to agentive decisions significantly increases in global supply chains where overlapping jurisdictions impose conditional and potentially conflicting requirements (Tong et al., 2022). As such, agentive systems must reason not only about what actions are permissible, but under which geographic, transactional, and temporal conditions those actions are permissible (Raji et al., 2020). Mapping regulatory requirements to agentive decisions therefore involves binding regulatory logic to state variables such as location, product classification, ownership structure, and transaction value (Chalendar et al., 2019). Such contextual binding enables agentive systems to dynamically assess their own compliance based on changes in their operational state (Hunt & Jackson, 2010), rather than relying on static rules.

Business Benefits of Effective Regulatory Mapping: Effective regulatory mapping provides businesses with reduced systemic exposure resulting from autonomous decision-making processes, which can accumulate over time through routine autonomous actions (Böhmecke Schwafert, 2024). Regulatory non-compliance often results from prolonged patterns of behavior rather than isolated incidents (Alles et al., 2008). If not for explicit mapping, autonomous systems may pursue strategies that may be operationally optimal but that are likely to exceed regulatory thresholds over time (Bose et al., 2014). Embedding regulatory logic within decision evaluation frameworks enables organizations to proactively prevent such drift and protect market access, revenue continuity, and corporate reputation (Böhmecke Schwafert, 2024).

Strategic Agility in Volatile Regulatory Environments: Additionally, regulatory mapping enables organizations to maintain strategic agility in environments characterized by frequent changes to regulatory requirements (Tong et al., 2022). Trade regimes, sanctions lists, and reporting standards are subject to change in response to geopolitical and economic events (Chalendar et al., 2019). Organizations that map regulatory requirements into decision architecture frameworks can easily update those mappings to accommodate changing regulatory requirements without having to redesign execution processes (Raji et al., 2020). As such, organizations with

strategically agile regulatory mappings can quickly adapt to changing regulatory requirements while continuing to execute autonomously (Hunt & Jackson, 2010). Agentive systems can therefore serve as tools of regulatory responsiveness rather than rigidities (Böhmecke Schwafert, 2024).

Transparency and Defensibility: Finally, regulatory mapping of agentive decisions supports organizational transparency and defensibility (Böhmecke Schwafert, 2024). When regulatory logic is explicitly tied to decision pathways, organizations can demonstrate how regulatory considerations affected decision-making outcomes (Jans et al., 2014). Such traceability is critical to organizational responses to regulatory inquiries, disputes, or audits (Alles et al., 2008). Autonomous decisions can therefore be explained not through rationalization subsequent to the event, but through documented regulatory assessment during execution time (Mitchell et al., 2019).

Theoretical Rigor in Regulatory Mapping: To achieve theoretical rigor in regulatory mapping, it is necessary to distinguish between regulatory prohibitions, conditional permissions, and reporting obligations (Raji et al., 2020). Some regulatory requirements prohibit specified actions (Tong et al., 2022). Other regulatory requirements permit actions under specified conditions (Chalendard et al., 2019). Finally, some regulatory requirements establish reporting obligations or documentation requirements without prohibiting execution (Alles et al., 2008). As such, regulatory mapping must classify regulatory requirements according to type and implement corresponding decision handling logic (Raji et al., 2020). This classification ensures that agentive systems respond appropriately to regulatory obligations (Böhmecke Schwafert, 2024).

Maintenance of Regulatory Mappings: Finally, regulatory mappings must be treated as a living document (Hunt & Jackson, 2010). Regulatory requirements evolve, interpretations of those requirements change, and enforcement priorities change (Tong et al., 2022). Governance frameworks must therefore support the ongoing updating and validation of regulatory mappings to ensure continued alignment (Vasarhelyi et al., 2004). Regulatory mappings that are static become outdated in dynamic regulatory environments (Bose et al., 2014). Ongoing maintenance of regulatory mappings therefore ensures that autonomous decision-making continues to align with evolving institutional expectations (Hunt & Jackson, 2010).

Compliance as a Decision Dimension: Ultimately, regulatory mapping of agentive decision-making into decision logic transforms compliance from an external constraint to an internal decision consideration (Raji et al., 2020). Regulatory considerations therefore become a factor in how agents evaluate decisions regarding actions, rather than an obstacle to action that occurs subsequent to decision-making (Böhmecke Schwafert, 2024). This integration of regulatory considerations into decision-making is critical to responsible deployment of agentive supply chains at scale in regulated environments (Böhmecke Schwafert, 2024).

Logistics and Trade Compliance

Trade Compliance and Logistics represent an operationally risky and legally significant element of agency-based Supply Chains; especially given that autonomous decision-making is made directly within the bounds of sovereign jurisdictions, customs regimes and international trade law (Chalendard et al., 2019). Traditional Supply Chain logistics compliance is monitored manually, by brokers and via post-shipment verification of shipping compliance (Voss & Williams, 2013). In traditional supply chains decision velocity is relatively slow, allowing for manual monitoring of compliance (Voss & Williams, 2013). However, in an agency-based supply chain all of the decisions regarding routing, carrier selection, consolidation, etc. occur continuously and at machine-speed (Alles et al., 2008). This continuous and rapid decision-making process results in a need for compliance logic to be integrated into autonomous decision-execution. Any delay in the evaluation of compliance will expose the organization to unacceptable levels of both financial and legal liability (Hunt & Jackson, 2010).

At its core, logistics and trade compliance involve ensuring that all movements of goods adhere to applicable trade regulations, customs requirements and transportation laws (Chalendard et al., 2019). Autonomous routing decisions must take into consideration various trade regulations including export controls, sanctions regimes, embargoed regions and preferential trade agreements that differ based upon the classification of products, the origin and destination of products and the ownership structure of the products (Tong et al., 2022). In an agency-

based system, compliance regulations cannot be used as an external validation after routing optimization (Raji et al., 2020). Rather, compliance constraints must be evaluated simultaneously with cost, time and capacity objectives (Hunt & Jackson, 2010). Integration of compliance with cost, time and capacity will ensure that no routes that violate trade law are ever considered feasible options, regardless of their operational appeal (Chalendard et al., 2019).

The complexity of trade compliance is further exacerbated due to the numerous layers of regulatory oversight present within global supply chains (Tong et al., 2022). Goods moving throughout a supply chain can pass through multiple countries, each requiring unique documentation and compliance measures (Chalendard et al., 2019). Autonomous logistics agents must evaluate this layered regulatory environment in real-time to incorporate customs clearance rules, bonded warehouse constraints and transit country regulations into their decision logic (Chalendard et al., 2019). Failing to do so may result in delayed shipments, fines, or even seizure of goods (Voss & Williams, 2013), resulting in an inability to realize any of the benefits derived from autonomous optimization (Voss & Williams, 2013). Embedding this decision logic into agency-based supply chains transforms compliance from an afterthought to a primary consideration when determining the feasibility of logistics (Tong et al., 2022).

Logistics compliance also includes transportation safety and labor regulations that determine how goods are moved versus simply where goods are moved (Tong et al., 2022). Autonomous agent selections and scheduling decisions regarding drivers and vehicles impact driver working hours, hazardous materials handling, vehicle certifications and route safety requirements (Tong et al., 2022). Many jurisdictions impose both civil and criminal liability for violations of transportation safety and labor regulations (Tong et al., 2022). Agency-based supply chains that solely focus on optimizing delivery speed and cost without evaluating safety constraints, run the risk of violating safety and labor regulations at a large scale (Böhmecke Schwafert, 2024). Embedded safety and labor compliance into autonomous decision logic will ensure that organizations maintain their commitment to their optimization objectives, while also upholding their legal and moral responsibilities (Raji et al., 2020).

Environmental compliance has grown significantly in importance for logistics governance as regulatory bodies mandate emissions reporting for transportation activities and implement carbon reduction initiatives (Sabeti et al., 2019). Autonomous routing and mode selection decisions directly influence the emissions profile of goods being transported by distance traveled, vehicle type and consolidation strategies (Sabeti et al., 2019). Within agency-based supply chains, environmental compliance must be evaluated in conjunction with logistics decision-making, rather than through post-hoc sustainability reports (Hunt & Jackson, 2010). Incorporating emission thresholds and reporting triggers into autonomous decision-making will ensure that agencies contribute to the organizational commitment to sustainability while complying with evolving environmental regulations (Sabeti et al., 2019).

From a business impact perspective, the consequences of failure to comply with logistics and trade regulations far exceed the consequences of the individual decision(s) made (Voss & Williams, 2013). A single shipment that fails to comply with trade regulations can lead to customs audits, increased inspection frequency, or suspension of trusted trader status, impacting future shipments within the same supply chain (Voss & Williams, 2013). In an agency-based system, where decisions are repeated thousands of times, failure to comply can quickly escalate into systemic disruption (Alles et al., 2008). Embedding compliance logic into autonomous execution will protect revenue continuity, preserve supplier and carrier relationships and safeguard organizational reputation within highly visible global markets (Tong et al., 2022).

In addition to regulatory requirements, logistics and trade compliance must also consider contractual obligations that overlap with regulatory requirements (Voss & Williams, 2013). Freight forwarding agreements, carrier contracts and service level agreements impose restrictions on routing, consolidation, subcontracting and liability allocation (Voss & Williams, 2013). Autonomous logistics agents must honor contractual obligations along with statutory regulations (Tong et al., 2022). Failure to comply with contractual obligations may not result in regulatory penalties but can result in litigation, loss of preferred pricing and/or termination of strategic partnerships (Voss & Williams, 2013). By integrating contractual compliance into agency-based decision logic, organizations can align their regulatory obligations with their commercial relationships (Raji et al., 2020).

The technical challenges of logistics and trade compliance arise from reconciling conflicting regulatory regimes (Tong et al., 2022). For example, actions permitted by trade regulation may violate environmental policies, or a compliant route may violate labor regulations in a transit country (Tong et al., 2022). Agency-based systems must reconcile these conflicts based upon defined organizational priorities and risk tolerance (Raji et al., 2020). To perform this reconciliation, agency-based systems require explicit governance logic that prioritizes compliance dimensions and defines acceptable trade-offs (Raji et al., 2020). Without such logic, autonomous systems may act erratically and unpredictably under complex regulatory conditions (Böhmecke Schwafert, 2024). Logistics and trade compliance in agency-based supply chains also require the creation of accurate and timely documentation related to the shipment (Chalendard et al., 2019). Documentation, such as customs declarations, certificates of origin, safety manifests and other documents related to transportation must be created accurately and consistently for each shipment (Chalendard et al., 2019). Therefore, agency-based execution systems must embed documentation workflows into decision pathways, ensuring that compliance-related documentation is created as part of the execution of the decision, rather than as an administrative task downstream of the decision (Hunt & Jackson, 2010). Integration of documentation into the execution of decisions will reduce errors in documentation, accelerate clearance and improve audit preparedness (Alles et al., 2008).

Strategically, embedding logistics and trade compliance into agency-based supply chains allows organizations to engage in global commerce with confidence while maintaining the ability to be agile in response to market opportunities (Tong et al., 2022). Organizations whose autonomous systems include compliance logic can rapidly respond to market opportunities without assuming excessive regulatory risk (Böhmecke Schwafert, 2024). This ability to respond to market opportunities transforms compliance from a perceived barrier to innovation into an enabler of scalable global operations (Böhmecke Schwafert, 2024). As such, agency-based supply chains can be faster, more efficient, and more robust and institutionally credible in complex regulatory environments (Böhmecke Schwafert, 2024).

Continuous Compliance Monitoring

Continuous Compliance Monitoring is a necessary element for any Supply Chain to operate in Regulated Environments (Vasarhelyi et al., 2004). Most Supply Chains are subject to various Regulations that are meant to protect consumers and the environment. In most cases these regulations require companies to take certain steps to comply with them. If a company fails to comply with these regulations it could result in fines, criminal charges, loss of licenses, etc.

Traditional auditing methods have been used for decades to ensure that companies are complying with regulations. These methods include periodic audits and inspections. Auditors will review documents and interview employees to determine if they are complying with all applicable regulations. However, these methods do not account for the fact that Supply Chains are constantly changing due to market fluctuations, global pandemics, natural disasters, etc. Therefore, traditional auditing methods may not be effective in detecting compliance failures.

The use of Continuous Compliance Monitoring allows companies to monitor compliance throughout the Supply Chain, in real-time. This includes the ability to monitor the movement of goods across borders, as well as the ability to detect changes in the composition of products being shipped (Alles et al., 2008). This type of technology uses AI to continuously scan and analyze the data flowing through the Supply Chain and make determinations of whether or not the data indicates that there is a potential compliance failure. If a potential compliance failure is detected, the AI will send a notification to the responsible employee(s) so that they can investigate the issue. Regulatory Risk is one of the biggest risks facing companies today. Many companies face large fines and other negative consequences when they fail to comply with regulations. However, traditional auditing methods may not be able to prevent compliance failures (Alles et al., 2008).

In addition to detecting compliance failures, Continuous Compliance Monitoring also helps to reduce the costs associated with compliance. Companies spend millions of dollars each year to audit and inspect their Supply Chains to ensure compliance with regulations. However, using Continuous Compliance Monitoring can help to

reduce the number of audits and inspections required, which can save companies money. Continuous Compliance Monitoring can also help companies to identify potential compliance failures before they occur, which can save companies even more money in the long run (Alles et al., 2008). In order to be effective, Continuous Compliance Monitoring must be designed carefully. If the monitoring is too broad, it can lead to false positives and unnecessary notifications. On the other hand, if the monitoring is too narrow, it may miss important compliance failures. Therefore, it is crucial to design the monitoring carefully to avoid either of these extremes (Alles et al., 2008).

Continuous Compliance Monitoring is also beneficial for improving customer satisfaction. Companies that use Continuous Compliance Monitoring tend to have higher customer satisfaction ratings compared to those that do not use Continuous Compliance Monitoring. There are several reasons why Continuous Compliance Monitoring improves customer satisfaction. First, customers expect companies to have high quality products that meet all relevant safety and environmental regulations. Second, companies that use Continuous Compliance Monitoring can respond faster to product recalls and other regulatory issues. Third, companies that use Continuous Compliance Monitoring tend to have fewer recalls and other regulatory issues compared to those that do not use Continuous Compliance Monitoring. All of these benefits improve customer satisfaction (Alles et al., 2008).

Continuous Compliance Monitoring also provides improved transparency and accountability. Since all transactions are monitored in real time, companies can prove compliance to regulatory agencies, suppliers and customers. This can help to increase trust among all stakeholders, and reduce the likelihood of disputes and lawsuits related to compliance (Alles et al., 2008). Continuous Compliance Monitoring also has the benefit of enabling adaptive governance in volatile regulatory environments. As trade regimes, sanctions lists, environmental standards and data protection rules evolve rapidly in response to geopolitical and economic conditions, traditional compliance rule-sets quickly become outdated (Tong et al., 2022). Continuous Compliance Monitoring Systems, however, enable organizations to continuously assess real-time behavior against up-to-date regulatory criteria, allowing for adjustments to autonomous execution to mitigate new regulatory exposure without disrupting ongoing operations (Hunt & Jackson, 2010). This adaptability is critical to maintain competitiveness and compliance.

Finally, the integration of Continuous Compliance Monitoring into agentic Supply Chains fundamentally alters how Compliance Functions collaborate with Operations (Vasarhelyi et al., 2004). Instead of serving as Gatekeepers approving/rejecting individual decisions, Compliance Teams now serve as Designers/Interpreters of Monitoring Frameworks defining Indicators Thresholds, and Aggregation Windows reflecting regulatory priorities and business risk tolerance (Vasarhelyi et al., 2004). Autonomous Systems subsequently enforce these definitions at scale (Alles et al., 2008). This transformation enables Compliance Experts to influence System Behavior Continuously, without becoming an operational Bottleneck (Hunt & Jackson, 2010).

However, Effective Continuous Compliance Monitoring is dependent upon Abstraction to avoid both Under-Sensitivity and Alert Fatigue (Hunt & Jackson, 2010). Monitoring Systems must therefore strike a Balance between Granularity and Interpretability, by aggregating Signals into Meaningful Indicators that Reflect Regulatory Risk (Vasarhelyi et al., 2004). Overly granular monitoring results in excessive Low-Level Alerts overwhelming Oversight Functions and eroding Trust in Autonomous Systems (Böhmecke Schwafert, 2024). Conversely, Under-Sensitive Monitoring enables Risk Accumulation without Detection (Alles et al., 2008). Thus, designing Appropriate Indicators depends on a Deep Understanding of Regulatory Intent, Supply Chain Dynamics and Business Impact (Vasarhelyi et al., 2004), requiring Governance Expertise Embedded within Technical Architectures (Böhmecke Schwafert, 2024).

Furthermore, Continuous Monitoring Supports Transparency and Accountability in Autonomous Supply Chains (Böhmecke Schwafert, 2024). Through Real-Time Visibility into Compliance Posture, Organizations can Demonstrate to Regulators, Partners, and Internal Stakeholders that Autonomous Execution is Actively Governed (Raji et al., 2020). Monitoring Dashboards and Reports Provide Evidence that Compliance is Not Assumed But Continuously Evaluated (Vasarhelyi et al., 2004). This Transparency Builds Institutional Confidence and Reduces Resistance to Autonomy Adoption (Böhmecke Schwafert, 2024). Moreover, it

Enhances Organization's Position During Regulatory Engagement by Demonstrating Proactive Risk Management (Böhmecke Schwafert, 2024).

The Analytical Foundation of Continuous Compliance Monitoring Can Be Formulated Using Composite Risk Evaluation Functions That Aggregate Multiple Compliance Indicators Over Time (Vasarhelyi et al., 2004). Specifically, let C_t Represent a Compliance Exposure Score at Time T Derived from Recent Decision Activity (Vasarhelyi et al., 2004). This Score Can be Represented as:

$$C_t = \sum_{i=1}^n w_i v_i(t)$$

Where v_i Represent Individual Compliance Indicators Such as Jurisdictional Exposure Transaction Volume or Emissions Contribution, and w_i Represent Their Relative Regulatory Importance (Hunt & Jackson, 2010). This Representation Emphasizes that Compliance Is Multi-Dimensional and Cumulative Rather Than Binary (Vasarhelyi et al., 2004). Additionally, It Enables Threshold-Based Governance Responses When Exposure Exceeds Acceptable Levels (Hunt & Jackson, 2010). Continuous Compliance Monitoring Also Enables Learning and Improvement Within Agentic Supply Chains (Vasarhelyi et al., 2004). Patterns Identified Through Monitoring Inform Refinement of Governance Constraints Policy Spaces and Escalation Thresholds (Hunt & Jackson, 2010). Over Time the System Becomes Better Aligned with Regulatory Realities and Business Objectives (Böhmecke Schwafert, 2024). This Feedback Loop Transforms Compliance From a Static Requirement Into a Source of Organizational Learning (Vasarhelyi et al., 2004). Autonomous Systems Evolve Not Only To Optimize Performance but to Internalize Regulatory Expectations More Effectively (Raji et al., 2020). In High-Risk Supply Chain Environments, Continuous Compliance Monitoring is Not Optional But Essential (Böhmecke Schwafert, 2024). Autonomous Execution Without Real-Time Regulatory Awareness Exposes Organizations to Unacceptable Legal and Reputational Risk (Böhmecke Schwafert, 2024). By Embedding Monitoring into the Decision Architecture, Agentic Supply Chains Achieve a Level of Regulatory Responsiveness That Manual Processes Cannot Match (Hunt & Jackson, 2010). This Capability Enables Organizations to Scale Autonomy Responsibly While Preserving Compliance Credibility and Institutional Legitimacy (Böhmecke Schwafert, 2024).

Audit Readiness

Audit readiness constitutes a critical basis for the institutional legitimacy of agentic supply chains, primarily due to the intersection of regulatory attention, contractually enforceable accountability, and public trust in these types of supply chains (Böhmecke Schwafert, 2024). Traditional supply chain auditing assumes that decisions are made by human entities who have express intent, appropriate judgment, and sufficient documentation so that those decisions can be retrospectively examined (Alles et al., 2008). However, the agentic nature of supply chains disassembles the assumptions of traditional auditing by employing autonomous systems that continuously make decisions in procurement, logistics, inventory management, and fulfillment (Raji et al., 2020). Therefore, the audit-readiness of agentic supply chains is concerned with the continued transparency, reconstructability, and defendability of autonomous decision-making processes during formal examinations (Böhmecke Schwafert, 2024). Consequently, the audit-readiness of agentic supply chains relies upon comprehensive decision traceability (Jans et al., 2014). Each autonomous action executed by an agentic system must be related to a verifiable record of the system's state at the time it was evaluated, the constraints employed, and the decision logic that was utilized to select each alternative (Mitchell et al., 2019).

Within supply chain operations, this could include the selection of sources, routes, inventory allocations, and compliance assessments (Raji et al., 2020). Decision traceability facilitates auditors' ability to reproduce decision-making paths after decisions are made, even though decisions occur in machine time and there has been no direct human involvement in the decision-making process (Jans et al., 2014). If traceability is absent, then autonomous decision-execution will become opaque, thereby diminishing regulatory confidence, and creating increased risk for non-compliance (Böhmecke Schwafert, 2024). Due to the velocity and scale of agentic supply chains, a systematic audit infrastructure is required to adequately support audit readiness (Alles et al., 2008). Autonomous systems can produce tens of thousands of decisions every hour in distributed networks across

multiple geographic locations (Alles et al., 2008). Manually reproducing activity of this scale is impractical (Jans et al., 2014). Therefore, audit-readiness is dependent upon automated logging architecture that captures execution events and governance evaluations in real-time (Vasarhelyi et al., 2004). Automated logging architectures must be designed to preserve temporal orderings, contextual information, and decision dependencies; otherwise, audit records will only represent that actions occurred, not how those actions arose from the decision framework (Rozinat & van der Aalst, 2008). Additionally, audit-readiness is dependent upon the alignment between declared governance policies and actual system behavior (Rozinat & van der Aalst, 2008). Auditors assess whether organizations adhere to their declared policies, regulatory mappings, and risk controls in practice (Alles et al., 2008). The alignment between declared governance policies and actual system behavior is realized in agentic supply chains via embedded governance and rule-constrained policy spaces, which ensure that execution behavior conforms to stated rules by construction (Raji et al., 2020). Stronger audit-readiness exists when execution logs can be directly related to governance definitions, and therefore, demonstrate that autonomous actions were permissible according to the governing framework applicable at the time (Mitchell et al., 2019).

From a business perspective, stronger audit-readiness minimizes the operational disruptions resulting from regulatory reviews and compliance investigations (Alles et al., 2008). Organizations that maintain well-structured audit data can respond to inquiry requests quickly, without incurring substantial managerial or operational costs (Vasarhelyi et al., 2004). This efficient response to inquiries reduces the indirect costs of compliance and permits supply chain operations to operate with minimal interruption (Alles et al., 2008). On the other hand, poor audit-readiness typically leads to protracted investigations, delayed shipments, and reputation damage that far outweigh the costs of proactive audit preparation (Böhmecke Schwafert, 2024). Audit-readiness is also pivotal in promoting contractual accountability among various stakeholders within complex supply networks (Voss & Williams, 2013). Most disputes regarding contractual obligations arise because of questions concerning whether parties complied with terms and conditions established for specific situations (Voss & Williams, 2013). Autonomous decision records provide tangible evidence of how decisions were made and whether terms and conditions were satisfied (Jans et al., 2014). As such, the evidence provided by decision records serves to reduce the potential for litigation and enhance an organization's negotiation position (Voss & Williams, 2013).

In agentic supply chains, audit-readiness therefore contributes to both regulatory compliance and commercial risk management (Voss & Williams, 2013). The technical requirements for audit-readiness extend beyond simple event logging (Vasarhelyi et al., 2004). Auditors generally require evidence of governance evaluations, which include the constraints employed, the alternatives evaluated, and why a particular action was chosen (Mitchell et al., 2019). Therefore, autonomous systems must simultaneously log governance decisions and execution actions (Raji et al., 2020). Dual logging provides auditors the ability to assess both the compliance of outcomes and the integrity of processes (Rozinat & van der Aalst, 2008). Absent governance-level evidence, auditors may determine that compliance was coincidental, rather than systematically maintained (Böhmecke Schwafert, 2024). Audit-readiness also enhances internal governance and organizational learning (Vasarhelyi et al., 2004). Periodic review of autonomous decision records enables organizations to recognize patterns of behavior that may necessitate governance refinement (Bose et al., 2014). Repeated borderline compliance decisions may indicate overly lenient constraints, or risk thresholds that are not aligned (Bose et al., 2014).

Audit analysis provides organizations with a feedback mechanism for improving governance frameworks continuously (Vasarhelyi et al., 2004). This learning function transforms audits from adverse events into opportunities for enhancing the strength of systems (Alles et al., 2008). Audit-readiness must accommodate diverse regulatory requirements and reporting standards in global supply chains (Tong et al., 2022). Different jurisdictions impose different documentation requirements, retention requirements, and audit formats (Chalendar et al., 2019). Therefore, agentic systems must provide flexible reporting capabilities to generate jurisdiction-specific audit artifacts that do not alter the fundamental decision-making logic (Mitchell et al., 2019). Flexibility is necessary to support the deployment of autonomous systems globally while maintaining uniform governance (Böhmecke Schwafert, 2024). Audit-readiness is essential for obtaining regulatory approval to deploy agentic systems (Böhmecke Schwafert, 2024). When organizations can demonstrate robust audit capabilities that ensure transparency and accountability, regulatory bodies are more likely to allow the use of

autonomous decision-making systems (Böhmecke Schwafert, 2024). Clear audit trails provide assurance to regulatory bodies that the use of autonomous systems does not preclude oversight (Raji et al., 2020).

In high-risk industries, such as pharmaceuticals, food distribution, or defense logistics, audit-readiness may be a requirement for deploying agentic systems (Böhmecke Schwafert, 2024). Audit-readiness also determines whether agentic supply chains can be sustained over time (Böhmecke Schwafert, 2024). While efficiency improvements gained through autonomy are fragile if they cannot survive regulatory scrutiny (Alles et al., 2008), embedding audit-readiness into system design ensures that autonomous decision-execution continues to be justifiable as regulations evolve, and enforcement intensifies (Böhmecke Schwafert, 2024). Ultimately, audit-readiness enables agentic supply chains to be sustained as technological innovations, rather than as transient technologies (Böhmecke Schwafert, 2024).

Auditability and Decision Traceability

Decision Provenance Logging

The ability to track the history of a decision (decision provenance) is fundamental to the auditing capabilities of an agentic supply chain; autonomous systems will make hundreds of decisions throughout the entire purchasing logistics inventory location and delivery cycle (Cheney et al., 2009). Auditing the provenance of a decision in a traditional supply chain is done via approval emails and procedural documentation accompanying each decision. However, in an agentic supply chain, since the decision-making is shifted to machine-based systems running faster and larger than human observation, the decision provenance has become the only method to maintain an organization's institutional knowledge of its autonomous actions so that decisions can be accounted for well after their execution (Cheney et al., 2009).

In terms of a supply chain, decision provenance goes beyond merely documenting the end result of the decision (such as routing choice or supplier assignment). Every autonomous decision in a supply chain is made based on evaluating a constantly changing and complex state of the system including but not limited to forecasting of future demand current inventory levels current capacity of suppliers contracted obligations regulations and/or other risk-related factors (Wang et al., 2022). Therefore, provenance must document the complete state of the system as the decision was being made to facilitate meaningful reconstruction of decision logic. If no documentation of the information environment exists in which a decision was made then there is no way for auditors and/or governance agencies to evaluate if an action was reasonable or compliant. Thus, provenance transforms decisions into richer artifacts that contain both the action and the context (Moreau et al., 2015).

As a direct result of the continuous nature of agentic execution, the need for structured provenance architectures is greatly amplified. Autonomous systems may produce tens-of-thousands of decisions per hour across many geographically dispersed nodes. The manual reconstruction of such activity is not feasible without the automation of provenance capture that is designed for high-volume and high-precision temporal resolution (Alongi et al., 2022). Provenance systems must also preserve the ordering relationship among decisions the temporal relationship between system state changes and decisions and the dependencies between/within agents. In supply chains, where decisions interact over time and across organizational boundaries the preservation of these relationships is essential for understanding downstream effects and systemic behaviors (Alongi et al., 2022). Provenance logging also plays a significant role in supporting regulatory compliance. As regulators increasingly require evidence that compliance considerations were evaluated when decisions were made versus after-the-fact, provenance records must include explicit references to compliance constraints regulatory checks and risk thresholds that were applied during the evaluation of a decision. Such evidence allows auditors to confirm that compliance was enforced systematically and continuously. In regulated supply chains provenance logging becomes proof of governance execution instead of a passive operational record (Wang et al., 2022).

Robust provenance logging also reduces the costs and disruptions associated with audits investigations and disputes. Companies with comprehensive provenance records can respond to regulatory inquiries quickly and efficiently without having to reconstruct decision logic retroactively. This capability minimizes operational disruptions and allows managers to continue focusing on core activities. Additionally, robust provenance records strengthen an organization's position in contractual disputes as they provide objective evidence of how and why

decisions were made under specific conditions (van der Aalst, 2016). Provenance logging also supports internal governance and performance management. Historical analyses of decision-making contexts provide insights into the types of governance gaps, optimization biases or unintended risk accumulations that may exist. For example, repeated decisions made under marginally acceptable compliance conditions may indicate overly permissive constraint settings or conflicting objectives. Governance bodies can develop more informed policies and threshold values using empirical evidence derived from provenance data. In this way, provenance logging can be viewed as a catalyst for continuous governance improvement (van der Aalst, 2016).

Designing provenance logging systems requires careful abstraction to ensure that provenance captures enough detail to be useful for auditing while maintaining the level of detail required for accountability. Logging all raw data signals would be impractical and unnecessary for auditing. Instead, provenance systems should capture only those state variables relevant to decision making, relevant governance evaluations, and the decision options considered. Selective capture ensures that logs remain interpretable to humans while containing sufficient detail for accountability. In supply chain operations, abstraction is especially important as excessive logging can obscure, rather than elucidate, decision rationale (Heluany et al., 2023). Additionally, decision provenance facilitates accountability in multi-agent supply chain environments where outcomes emerge from interactions between autonomous decisions. By linking decisions to the specific agents responsible for scope and context of each decision provenance logs provide organizations with the means to assign accountability within the system. Accountability assignments are essential to governance oversight as they differentiate between the local decision behavior and the overall system effect. Accountability assignments also allow for targeted corrective actions and avoid blanket restrictions on autonomy that could negatively affect system performance (Wu et al., 2022). In global supply chains, provenance logging must also consider jurisdiction-specific requirements related to data retention, access, and disclosure. Various regulatory jurisdictions have different requirements for how long decision records must be retained, who can access them, and what can be disclosed. Therefore, provenance architectures must provide configurable retention policies while preserving data integrity. This flexibility allows organizations to meet a variety of regulatory obligations without partitioning their decision-making systems (Heluany et al., 2023).

Provenance logging also increases the willingness of managers and stakeholders to delegate authority to agentic systems. When managers and stakeholders know that decisions made by autonomous systems can be reconstructed and reviewed if needed, they are more likely to grant authority to such systems. Provenance provides this assurance by ensuring that autonomy does not equate to opacity; decisions made by machines remain observable and reviewable within institutional governance frameworks (Kuehn, 2018). To summarize, decision provenance logging transforms autonomous decision-making into a transparent, auditable, and accountable organizational process, rather than an ephemeral algorithmic activity. The tracking of the lineage context, and governance evaluations of each decision, through provenance logging provides assurance that agentic supply chains are compliant with regulatory obligations. Provenance logging is a foundation upon which the trust of stakeholders, and regulatory acceptance of autonomous decision systems in complex supply chain environments, can be built (Kuehn, 2018).

Explainable Policy Execution

The need for explainable policy execution exists so that organizational stakeholders can understand how an organization's autonomous decision-making processes relate to its objectives and governance structures. Autonomous systems make operational decisions — e.g. sourcing allocation, routing, and fulfillment prioritization — that affect the operation of the supply chain. Organizational stakeholders — including managers, auditors, and regulators — require clear and understandable explanations regarding how these decisions relate to business practices and governance requirements. Explainability bridges the gap between algorithmic decision-making and institutional decision-making by explaining how autonomous decisions are made based on the organization's objectives and constraints (Ribeiro et al., 2016).

Explainability of autonomous decisions must be articulated in terms of business practices and not computational abstractions. Organizational stakeholders do not care about internal model parameters; however, they are interested in knowing if the autonomous decisions made by the system respect established business practices —

i.e. did the decision respect cost constraints? Did the decision respect service commitments? Did the decision respect risk thresholds? Were the decisions compliant with applicable laws and regulations? Explainable policy execution addresses the above questions by providing explanations of what criteria were used to make each decision and how trade-offs were addressed. An example of explainable policy execution would include explaining a sourcing decision through the lens of a supplier's reliability assessment; a supplier's eligibility under regulatory requirements; the costs associated with the suppliers; and the timeline for delivery.

Explainable policy execution also supports regulatory engagement. Regulatory bodies are requiring greater transparency into how decisions are made using automated decision-making systems to ensure compliance with applicable laws and regulations. In supply chains, this would include demonstrating that routing decisions respected trade restrictions; demonstrating that sourcing decisions complied with all applicable labor regulations; and demonstrating that allocation strategies complied with all applicable contracts. Explainable policy execution provides organizations with a means to demonstrate their compliance reasoning with their respective regulatory bodies without having to expose proprietary algorithms. This balance preserves competitive advantages of the organization while meeting regulatory scrutiny (Verma et al., 2024).

Explainability also has a positive impact on business adoption of autonomous systems. Supply chain professionals are more likely to trust and use agentic systems when they can understand and validate decision rationale. Explainability reduces concerns that autonomous systems will make decisions based on "hidden" objectives or that they will optimize narrowly at the expense of broader organizational goals. Trust increases the rate of deployment across critical supply chain functions and it promotes acceptance of autonomous decision-making throughout an organization (Ribeiro et al., 2016).

Explainable policy execution also enables the effective human oversight of agentic systems within governance frameworks. Supervisors who are responsible for overseeing agentic systems utilize explanations to evaluate any deviations in performance from anticipated behavior. The absence of explainability could result in over-intervention or disabling of autonomy due to false assumptions regarding deviations in performance. However, with explainability, supervisors can differentiate between legitimate adaptations to changes in the environment and legitimate governance drift. This differentiation improves the quality of supervision and preserves operational efficiency (Verma et al., 2024).

The theoretical basis for explainable execution requires the ability to distinguish between local explanations of individual decisions and global explanations of policy behavior. Local explanations provide justification for why a specific action was taken under specific circumstances. Global explanations provide an overview of how a decision policy behaves under a variety of circumstances. Both levels of explanations are required to determine whether an autonomous system is aligned with strategic objectives in a supply chain governance context. Therefore, effective autonomous systems should support multi-level explanations to meet varying oversight needs (Guidotti et al., 2018).

Explainability also facilitates dispute resolution in supply chain operations. Disputes arise between partners, customers, and regulators, when they question autonomous decisions. Explanation of decision rationale provides evidence of fairness, compliance, and contractual adherence. Therefore, explainability provides protection against potential litigation and damage to reputation. In cases where autonomous decisions have significant financial implications, explainability provides a safety net to protect against litigation and damage to reputation (Ribeiro et al., 2016).

Finally, scalability is a critical aspect of explainable policy execution. Large numbers of autonomous decisions are generated in agentic supply chains and manual explanation of decisions is impractical. The use of automated explanation generation using standard templates and governance criteria, ensures that explanations are consistent and efficient. Explanations can be stored along with provenance logs to facilitate auditing and oversight functions without adding operational burden (Verma et al., 2024).

Explainable execution also supports organizational learning. Reviewing explanation patterns enable organizations to recognize implicit biases or unintended trade-offs in the decision logic. Organizations can refine

objectives, constraints, and governance parameters based on this insight. Ultimately, explainability provides a means to continuously align autonomous behaviors with changing business priorities (Guidotti et al., 2018).

Therefore, explainable policy execution provides a means to transform autonomous decisions into institutional intelligible actions. By articulating decision rationale in terms of business practices, agentic supply chains become transparent, accountable, and governable. This transformation is essential for sustained adoption of autonomy in complex, regulated supply chain environments (Verma et al., 2024).

Replayable Decision Paths via Digital Twins

Agile decision pathways through digital twins, allow for an experiential component to auditability in agentic supply chains (Ivanov & Dolgui, 2021). Provenance logs and explanations provide static representations of decision logic, while digital twins provide the ability to simulate autonomous behavior in a simulated operational environment. The simulation of autonomous behavior in a simulated operational environment is important in supply chains, because autonomous decisions interact in non-linear ways across time-space and across different levels of organization. Digital twin replayability provides the opportunity to observe how an autonomous system responds to certain conditions instead of having to infer based on record history (Kritzinger et al., 2018).

Digital twins include the structural and behavioral attributes of physical supply chain characteristics such as inventory flows, transportation networks, suppliers' constraints and regulatory requirements (Piancastelli & Tucci, 2020). They provide the ability to synchronize with historical execution data to represent the state of the system at the time decisions were made. Therefore, digital twins enable auditors and analysts to play back the decision sequences as they occurred to gain deeper insight into decision dynamics and system behavior (Wang et al., 2022).

In particular, the replayable nature of decision paths are useful in understanding emergent outcomes in complex supply chains. Emergent outcomes in supply chains such as delays, shortages or excessive costs occur due to the interactions of many decisions versus one singular action. Digital twin replayability provides analysts the ability to observe how a sequence of autonomous decisions propagates throughout the network over time. A temporal view of this propagation is necessary to identify systemic weaknesses and governance gaps which static views may fail to find (Ivanov, 2020).

From a business viewpoint, digital twin replayability decreases friction during audits and investigations. Auditors are able to interactively examine the decision-making process of an autonomous system as opposed to being limited to documentation. This increased transparency will accelerate the auditing and investigation processes and increase confidence in autonomous systems. Additionally, digital twin replayability will reduce the burden on operational teams by allowing them to visually demonstrate the logic and outcomes of autonomous decision-making (Abouelrous et al., 2023).

Additionally, digital twin replayability enables counterfactual analysis. Organizations are able to modify inputs, constraints or policies within the digital twin to determine how alternative conditions would have resulted in different autonomous decisions. This is important to organizations in determining whether compliant alternatives existed, or whether decisions made were reasonable under the circumstances. Counterfactual replayability will strengthen the defense position of an organization in regulatory reviews and disputes (Verma et al., 2024).

Furthermore, digital twin replayability will facilitate governance validation and system testing prior to deployment of new policy or expansion of autonomy. Organizations will be able to test the effects of proposed changes to their policies, processes or constraints using historical scenarios prior to implementation. This testing will expose potential unintended consequences and ensure that governance and/or process changes are effective. Thus, digital twin replayability will decrease the risks associated with deployments and increase the overall reliability of a system (Burgos & Ivanov, 2021).

The accuracy of replayable decision paths depend upon maintaining a high-fidelity synchronization of the physical operation states with the digital representations of those states. Synchronization errors or lags undermine the validity of replay analysis. High fidelity requires robust architectures for data ingestion and state

management. This highlights the importance of digital twins as operational tools as opposed to illustrative tools (Wu et al., 2022).

Replayability will contribute to organizational learning and training. Governance and management personnel will be able to understand the dynamics of autonomous behavior through observing that behavior in realistic scenarios. The experiential learning derived from the observation of autonomous behavior will foster both trust and competency in overseeing autonomous systems (Kritzinger et al., 2018). Replayable decision paths will increase accountability by rendering autonomous behavior observable rather than abstract. Decision makers will be able to see how decisions were made and how they impacted outcomes. The observability provided by replayable decision paths will increase the legitimacy of institutions and will be key to long-term acceptance of autonomous decision-making (Ivanov & Dolgui, 2021). Therefore, digital twin replayability will elevate auditability from static inspections to dynamic verifications. Agentive supply chains will be not only traceable but experientially understandable, thereby enabling robust governance in complex environments (Kuehn, 2018).

Causal Attribution

Attribution of cause is important to the management of agentic supply chains because virtually all operational events in such systems are the result of the interactions of several autonomous decision-making entities. As a result, service failures, price increases, inventory imbalance, regulatory violations and downstream disruptions result from the sequence of decisions taken by each entity in response to uncertainty and partial visibility. In agentic systems there are many potential simultaneous influences on procurement, routing, allocation and fulfillment, and exogenous events such as port closures, supply interruptions and changes in demand create changing circumstances for entities. Attribution creates a defensible connection between decisions and their consequences, creating accountability based upon supply chain dynamics, rather than mere association (Ivanov, 2020). Such accountability is critical to organizational learning, regulatory credibility and risk governance in high consequence environments.

Attribution in supply chain operations must recognize that each decision is an intervention that modifies the path of the system's state at different points in time (Brodersen et al., 2015). An allocation decision will change the location of inventory and affect the feasibility of service in subsequent periods. A routing decision will change the distribution of lead times and will pass variability to production scheduling. A sourcing decision will increase the entity's vulnerability to suppliers and will change the probability of disruptions. Attribution must therefore include the modeling of how each decision is propagated through inventories, capacity and constraint, rather than simply whether each decision occurred at the same time as the event. Attribution includes a temporal and structural understanding of the propagation of decisions through a supply chain and is consistent with the notion that supply chains are coupled dynamic systems characterized by feedback loops and time delays (Ivanov & Dolgui, 2021).

Causal attribution has value to businesses in the form of the ability to differentiate between the occurrence of an event due to excessive risk realization versus failure of policy. A disruption may occur even if all decisions were prudent and compliant, especially in the presence of extreme uncertainty. Conversely, a large loss may occur due to the fact that the policy decisions implicitly concentrated risk, amplified volatility and ignored constraints in the edge cases. Accountability enabled by attribution allows governance bodies to identify whether the undesirable event was caused by an external shock that could not have been reasonably controlled, or by a policy decision that should be corrected. This differentiation is necessary to maintain trust in autonomy without expecting unachievable perfection (Burgos & Ivanov, 2021).

Causal attribution also provides the analytical framework for corrective action. Without attribution, organizations will generally react to failure by restricting autonomy broadly, or by granting blanket approval that will degrade performance. Attribution will allow for the identification of the need for correction of specific policy decisions, including the tightening of constraints for specific decision classes, the adjustment of escalation thresholds in specific contexts, or the modification of reward structures that encourage risky behavior. In supply chains where the business value of autonomy is dependent on speed and scale, targeted remediation will preserve

value while increasing safety and compliance. Attribution therefore serves as a stabilization mechanism for governance to prevent over-reaction and maintain institutional confidence (Brodersen et al., 2015).

Another advantage of causal attribution is its role in the auditability and defensibility of an organization's actions under regulatory and contractual review. Regulators and contract partners often do not just want to know what occurred, they want to know why it occurred and what decisions contributed to the materiality of the event. Attribution enables organizations to provide evidence-based narratives that link autonomous decisions to observed outcomes through traceable mechanisms. The use of attribution will decrease the possibility of adverse regulatory interpretation, strengthen positions in contractual disputes and provide a transparent report to boards and oversight functions. In regulated supply chains, causal attribution can serve as a requirement for increased levels of autonomy (Cheney et al., 2009). Causal attribution in agentic supply chains must also consider the interaction effects of multiple agents. Outcomes may be influenced by coordination failures that result in locally optimal decisions being combined to generate globally suboptimal outcomes such as oscillating inventory exchanges, congestion amplifications or cascading supplier switching events. Attribution must account for these interaction effects by examining the decision sequences and joint action patterns that generated the outcome, rather than attributing the outcome to an individual agent. This type of examination will support governance by determining whether problems result from flawed policy decisions or coordination architectures that require design revision. In complex networks, attribution that accurately identifies problem origins will enhance both technical improvement and organizational accountability (Abideen et al., 2021).

The formalization of causal attribution can be defined through a counterfactual marginal contribution of an action to an outcome. Let O denote an outcome of interest, such as total cost, delay, shortage rate, or exposure to compliance; and let a_1, a_2, \dots, a_n denote the realized sequence of decisions within a given horizon. The marginal causal contribution of decision a_k to an outcome O can be defined by comparing the realized outcome to a counterfactual outcome in which a_k has been substituted with a null or baseline alternative, while a_k is held constant.

$$\Delta O_k = O(a_1, \dots, a_k, \dots, a_n) - O(a_1, \dots, a_{k-1}, \emptyset, a_{k+1}, \dots, a_n)$$

Attribution defines the incremental difference between the realized trajectory and the counterfactual trajectory in which the decision is absent or neutralized (Verma et al., 2024). In supply chain contexts, the counterfactual term is generally calculated using a replay mechanism such as a digital twin that replicates system dynamics under the altered decision sequence. Attribution provides a governance function because it distinguishes between decisions that produced a material change in outcomes and decisions that were simply incidental to an outcome, thereby providing accountable responsibility, rather than general accountability (Brodersen et al., 2015).

Causal attribution also supports strategic performance optimization by demonstrating which decision classes produce the highest marginal impact on key outcomes. If decisions regarding routing produce consistently high marginal contributions to delay variance, governance can focus on policy improvements in transportation. If decisions related to switching suppliers produce consistently high marginal contributions to cost volatility, governance can revise the constraints on switching suppliers, or implement hysteresis mechanisms to reduce the frequency of switches. By identifying which decision classes have the greatest marginal impact on performance, attribution transforms attribution into a value instrument for business strategy, by directing investments towards the decision levers that most directly affect financial and service performance. The overall effect is a more efficient and stable autonomy regime (Guo et al., 2025).

Causal attribution must also account for the existence of uncertainty and partial observability. Supply chain outcomes represent stochastic disturbances and measurement errors, and counterfactual analysis is susceptible to model error. Therefore, governance frameworks consider attribution as probabilistic evidence rather than absolute certainty, combining attribution results with confidence measures based upon model fit, data quality, and variance of scenarios. This disciplined approach to accounting for uncertainty will improve the credibility of attribution results with auditors and regulators, as it avoids overstating deterministic causality while providing a structured accountability framework. In mature agentic supply chains, attribution results will be used in ongoing governance reviews and risk committee meetings (Ivanov, 2020).

Causal attribution provides the final layer of the auditability and traceability stack by establishing a causal relationship between the origin of decisions and the resulting outcomes through a defensible counterfactual argument (Brodersen et al., 2015). Provenance describes the history of the decisions made and the constraints under which the decisions were made; explainability describes the process by which the policies selected the actions taken; and attribution describes how the actions selected impacted the realized operational trajectory. Collectively, these three concepts provide a transparent and defensible autonomy that can be employed in regulated and high-risk environments. Causal attribution addresses a significant adoption barrier by providing a credible basis for accountability, remediation and continuous improvement of autonomous decision-making regimes (Brodersen et al., 2015).

Data Sovereignty in Global Agentic Supply Chains

Data Localization Constraints

Data localization rules now shape structural conditions for global supply chains, as governments establish authority over how data created in each country is stored, processed, and transmitted (Taylor, 2020). In agentic supply chains, which make decisions autonomously using real-time data, the ability to store process and transfer data is restricted due to data localization rules, resulting in architectural changes to the way intelligence is used and exercised in autonomous global supply networks.

As it relates to supply chain operations, the data required to be stored locally to meet localization requirements can include: transactional data; contractual agreements with suppliers; current inventory levels; shipment tracking data; and increasingly, sensor data and operational telemetry. The majority of countries require companies to store data either within its borders or process data only via equipment that is located within a country's boundaries (Taylor, 2020). Therefore, for agentic systems that utilize total visibility throughout their networks, data fragmentation limits the total information available to agents for decision-making. Agents, therefore, must operate under "partial observability" where some data cannot be aggregated centrally nor freely shared (Bernstein et al., 2002). The condition of partial observability changes the design architecture of decision intelligence for agentic systems.

Regulatory rationales for data localization often extend beyond the domain of privacy to include economic and strategic rationales (Hummel et al., 2021). Governments view supply chain data as a strategic asset that discloses industrial production capabilities, trade dependency relationships, and vulnerabilities in critical infrastructure. Therefore, governments impose restrictions to prevent foreign entities from accessing or controlling the data. If agentic supply chains fail to comply with localization requirements, they will face potential regulatory penalties, mandatory divestitures, or exclusion from major markets. Thus, compliance with localization requirements is not just a legal obligation, but a necessary condition for agentic supply chains to sustainably operate globally.

From a business perspective, data localization creates tradeoffs between the extent to which companies can exercise intelligence, and their ability to comply with regulatory requirements. Companies attempting to circumvent localization rules by transferring data informally, risk facing severe penalties and reputational damage. Conversely, if companies choose to overly segregate data to avoid regulatory issues, they risk degrading their ability to execute timely and responsive decisions. Therefore, agentic systems must be designed to respect the limitations imposed by localization rules while maximizing the use of locally available data (Kairouz et al., 2021). Achieving this balance will require the development of sophisticated architectures versus relying solely on policy-level enforcement.

Data localization rules also impact the governance and accountability of autonomous systems. Agents' decisions must be traceable back to the data sources that exist within a jurisdiction's legal framework. Provenance records must provide evidence of not only how decisions were made, but also where the underlying data existed at the time of the decision. Thus, data localization rules expand auditability to include the spatial dimensions of data governance. As such, autonomous systems must maintain knowledge of the location of the data (residency) as part of the decision-making context (Hummel et al., 2021). The relationship between data localization constraints and real-time decision-making poses additional challenges in logistics and fulfillment. Decisions regarding routing and allocation typically involve cross-regional coordination under stringent time constraints. Data

localization rules can limit access to upstream or downstream data that exists in another jurisdiction. Therefore, agentic systems must develop alternative solutions including predictive modeling, local approximations, or delayed synchronization to mitigate the effects of data localization. Such workarounds add uncertainty to the decision-making process, which must be addressed through explicit representations in decision logic (Zhang et al., 2022).

Data localization also impacts the strategies of vendors and platforms. Cloud-based centralized architectures may be infeasible in jurisdictions where sovereign infrastructure is required. Therefore, organizations may be required to implement regional compute clusters that are operated by local authorities (Bonomi et al., 2012). Therefore, data localization rules force agentic supply chains to pursue decentralized infrastructures as a governance requirement versus a performance optimization. Such decentralization has implications for costs and complexity that must be integrated into a company's business strategy. Furthermore, data localization rules create additional complications in responding to incidents and managing risks. When companies experience disruptions, they may seek to collect data quickly and efficiently from across multiple regions to support their response efforts. However, data localization rules may prevent the aggregation of data from different regions during times when companies require global visibility (Taylor, 2020). Therefore, agentic systems must be designed with built-in mechanisms for accessing data during emergencies while complying with data localization rules, or with local autonomous response that does not rely on a centralized control structure (Taylor, 2020).

In the long run, data localization rules will lead to a transition from monolithic global intelligence to regionally autonomous decision ecosystems. Therefore, agentic supply chains will evolve into a federation of local intelligence nodes that share data in a constrained manner (Yang et al., 2019). This evolution will challenge traditional notions of global optimization, but aligns the autonomy of agentic systems with the geopolitical realities of localized data storage (Hummel et al., 2021). Ultimately, data localization rules define the boundary between what agentic global supply chains can know versus what they must infer. Instead of viewing localization as an afterthought, organizations should treat localization as a first-class architectural constraint to enable them to build autonomous systems that are both compliant and resilient. Building such autonomous systems is critical to sustaining agentic supply chains in an environment of digital sovereignty (Hummel et al., 2021).

Sovereignty Aware Digital Twin Partitioning

Sovereignty-aware digital twin partitioning takes data localization ideas further in the modeling and simulation layers of agentic supply chains (Tao, Zhang, Liu, & Nee, 2019). Sovereignty-aware digital twins traditionally have been assumed to be part of a unified supply network representation in which all related data are consolidated into a single model. However, in the case of sovereign entities, this assumption will no longer hold. Digital twins must therefore be partitioned to recognize the legal and jurisdictional boundaries of each entity while still enabling autonomous decision-making. Thus, partitioning will become a fundamental design feature of digital twins rather than an implementation feature. As supply chain applications digital twins model inventory flows production capacity transportation networks, and contractual relationships (Tao, Cheng, Qi, Zhang, Zhang, & Sui, 2018), sovereignty-aware partitioning must segment these models based on the rules governing data residency. Each regional twin will run on data permitted by law in its respective region while either abstracting or anonymizing dependency on other regions. Segmentation will provide assurance that both simulation and decision validations occur within the bounds of applicable laws, while simultaneously providing operational relevance.

The way in which scenario analysis and stress tests are performed will also change with partitioned digital twins. Global scenarios will be broken down into regional simulations which will communicate through constrained interfaces. For example, a disruption in one country may be represented to another region as an aggregate signal representing the potential impact of that event, rather than as raw data. Therefore, agentic systems will be required to use abstraction techniques to understand cross-regional impacts without having direct knowledge of those impacts. The importance of designing strong interfaces between twin partitions will thus increase (Tao, Zhang, Liu, & Nee, 2019). From a regulatory point of view, partitioning provides sovereignty-aware digital twins greater regulatory defensibility. Organizations can demonstrate that they do not send their sensitive operational data outside of their own jurisdiction, while at the same time maintaining the capability to make

decisions regarding that data. Transparency regarding the handling of operational data will reduce regulatory barriers, build confidence in organizations' commitment to compliance, and create trust with authorities. Digital twins will therefore function as compliance instruments as well as operational tools (Tao, Cheng, Qi, Zhang, Zhang, & Sui, 2018).

Partitioning also will affect the workflows used for validating decisions. Policies developed or proposed in one region cannot be validated using the full global state; instead, validation must occur using the local twin partitions along with synthetic representations of external behavior. This constraint means that care must be taken in developing validation metrics that will assure that decisions made under incomplete information remain valid. Any governance framework applied to these digital twins must take this constraint into consideration. Business-wise, the partitioning of digital twins may result in additional costs associated with infrastructure and coordination complexity. Multiple regional twins must be operated, synchronized and managed/governed. However, the benefit derived from being able to operate in jurisdictions that could not allow advanced analytics or autonomy, will likely outweigh the added expense. In essence, the partitioning of digital twins will open up markets to organizations that could not previously operate in those markets due to restrictions on advanced analytics or autonomy.

Partitioning also adds to the resiliency of agentic supply chains. Even if cross-border connectivity is severely limited or eliminated, regional twins can continue to operate autonomously. This capability is especially beneficial in geopolitical risk scenarios where data flows may be intentionally or unintentionally disrupted. As such, agentic supply chains that have partitioned digital twins will be more resilient during severe disruptions (Bonomi et al., 2012).

However, partitioning does present some challenges to the learning and optimization functions. Global patterns may be more difficult to identify and utilize when data is partitioned. Agentic systems will therefore have to rely on federated insights or shared abstractions rather than pooling raw data. The limitations presented here reinforce the necessity of developing sophisticated mechanisms for coordinating activities across twin partitions (Zhang et al., 2022). Finally, from an organizational viewpoint, sovereignty-aware digital twin partitioning will significantly change how global supply chain governance is practiced. Control will be decentralized and distributed among multiple regional autonomy entities. This trend in control is consistent with broader trends in geopolitical decentralization and regulatory fragmentation. Therefore, sovereignty-aware partitioning transforms digital twins from monolithic mirrors into modular governance-aware simulation environments. The transformation is necessary to accommodate agentic supply chains and their autonomous decision-making capabilities within the confines of data sovereignty realities (Tao, Zhang, Liu, & Nee, 2019).

Federated Learning

Agencies face an inherent trade-off when optimizing their supply chains across borders. The typical classical method of optimization, referred to as global optimization, considers the entire network, coordinates all decision-making functions throughout the network, and shares all information regarding the network. However, agencies do not have full control over their networks; they are limited by data sovereignty requirements, which prohibit them from sharing data across borders and make it impossible for them to have centralized control over their networks. As such, agencies must consider trade-offs between achieving locally optimum solutions using the available information and achieving globally optimum solutions (Bernstein et al., 2002). Supply chains typically use local optimization for supply chain operations. While local optimization can enhance an agency's ability to respond to changes within its region of responsibility and increase the agency's compliance with regulations at the local level, it may decrease the global efficiency of the supply chain. For example, the use of regional inventory optimization could improve the agency's ability to meet local demand while simultaneously increasing the amount of inventory held by the agency and thus decreasing the overall efficiency of the supply chain.

On the other hand, a global optimization solution would likely result in decreased costs for the supply chain, but would most likely violate the agency's localization and regulatory requirements. Thus, agencies must take into consideration both the local objectives of each segment of the supply chain and the global objectives of the

supply chain as a whole. In particular, agencies should consider how their local optimization objectives will impact their global optimization objectives.

From a theoretical perspective, the trade-off between local optimization and global optimization can be framed as a trade-off between the optimization space defined by the agency's data sovereignty requirements and the global optimization space. Clearly, no global optimum can be achieved if data cannot be shared freely across the organization. As such, the agency must find a constrained optimum that respects the agency's jurisdictional boundaries. The shift in the optimization problem from unconstrained global efficiency to compliant network performance (Hummel et al., 2021), highlights the need for the agency to consider the trade-off between local and global objectives as part of the optimization process. Failure to manage this trade-off effectively can lead to serious consequences for both the compliance and performance of an agency. On one hand, failure to coordinate decisions among segments of the supply chain can lead to inefficient decision making and a lack of economies of scale. On the other hand, failure to comply with data sovereignty requirements can lead to regulatory sanctions, loss of market access and damage to reputation. Agencies must therefore develop a set of logical rules that capture the trade-off between local and global objectives and reflect their organization's priorities, risk tolerance and regulatory position.

As discussed above, digital twins and federated learning offer the potential to help navigate the trade-off between local and global objectives by allowing agencies to achieve approximately global reasoning and coordination without having to share data directly. Through the use of digital twins and federated learning, local agents can optimize based on local state while coordinating through abstract signals or shared models (Irfan et al., 2024). The coordination enabled through digital twins and federated learning allow local agents to achieve some degree of alignment with the global objectives of the supply chain while respecting the agency's sovereignty constraints.

The trade-off between local and global objectives can be formally expressed as a constrained optimization objective. A global performance metric J_g represents the agency's global objectives while a vector of regional performance metrics J_r represents the regional objectives of each segment of the supply chain. Sovereignty constraints $S_r, r \in R$ define the agency's data sovereignty requirements. The constrained optimization objective is:

$$\max_{\pi} J_g(\pi) \text{ subject to } \pi \in \bigcap_r S_r$$

The constrained optimization objective highlights that the agency's global objectives can only be achieved within the intersection of the agency's regional sovereignty constraints. The constrained optimization objective makes it clear that feasible policies are defined by jurisdictional rules rather than purely technical considerations (Yang et al., 2019). The trade-off between local and global objectives has implications for governance and accountability. Decisions made by regional units may appear to be reasonable from a local perspective, but ultimately lead to suboptimal global results. To address this issue, governance frameworks must establish acceptable levels of global inefficiency in exchange for compliance and resilience. Establishing such standards must be done in a transparent manner to avoid ambiguity regarding accountability (Taylor, 2020). Finally, the trade-off between local and global objectives has implications for the organizational structure of an agency. Decision authority can be devolved to regional units with coordination mechanisms instead of relying on centralized command. An organizational structure that aligns autonomy with sovereignty requires effective governance to mitigate the risk of fragmentation (McMahan et al., 2017).

From a resilience perspective, regional autonomy provides greater robustness against disruptions caused by geopolitical events. If the agency is unable to coordinate decisions across borders due to disruptions, regional units can continue to operate independently. The resilience benefits of regional autonomy must be considered relative to the efficiency losses associated with decentralized decision making during normal times (Zhang et al., 2022). In summary, the trade-off between local and global optimization defines the strategic boundaries of autonomous supply chains under sovereignty constraints. Instead of considering the trade-off as an implicit constraint of optimization, agencies can treat the trade-off as a design variable to create autonomous systems that are both compliant and competitive (Haripriya et al., 2025).

Risk Containment and Failure Management

Agent Drift Detection

Agent drift detection has been identified as a fundamental function for maintaining consistency among autonomous decision-making processes and organizational goals in agentic supply chains (Gama et al., 2014). Agentic systems dynamically update their internal decision-making rules based on changing data distributions, operational feedback, and environment volatility (Bifet & Gavalda, 2007), which is required for performance in supply chains subject to fluctuating demand; geopolitical uncertainty; and supplier uncertainty. However, the adaptive nature of these systems poses a risk that their decision-making processes will diverge from business objectives; compliance requirements; and/or risk tolerance over time. Drift detection serves as a means to ensure that learning continues to be consistent with governance instead of uncontrolled.

In operational supply chain settings, initial symptoms of drift may manifest themselves in gradual changes in priorities rather than catastrophic failure. For example, an agent might slightly favor less expensive transportation routes resulting in marginal increases in regulatory exposure or delivery variability. Alternatively, another agent may determine that meeting customer service requirements is paramount at the expense of inventory efficiency which results in long-term cost inflation. As noted above, many of these types of deviations typically remain undetectable using aggregated performance metrics until these deviations are manifested in operational or compliance failures. Therefore, the primary focus of drift detection is to monitor behavioral trends as opposed to evaluating singular performance events.

To effectively identify drift in operational supply chain settings, it is imperative to define dynamic boundaries of acceptable behavior as opposed to rigid baselines. Because supply chain environments are inherently non-stationary, fixed baselines rapidly become outdated. Drift detection systems must assess whether observed behavior continues to fall within the context-adjusted boundaries that have been established through governance objectives and operational conditions (Truong et al., 2020). The use of dynamic boundaries provides a mechanism to distinguish healthy adaptation from unhealthy deviation allowing for the preservation of learning benefits while mitigating misalignment.

From a governance standpoint, drift detection represents a paradigmatic shift from episodic oversight to continuous assurance. Traditionally, governance is conducted through periodic audits or performance reviews that are insufficient to address the rapid decision-making pace of machines. Instead, drift detection embeds oversight directly into the system architecture providing early warning signals as behavior changes. Therefore, drift detection is consistent with governance by design principles and facilitates proactive remediation. The potential for business-related impacts due to drift detection is considerable since it can prevent slow-moving failures that are difficult and costly to correct. A number of supply chain disruptions associated with automation have occurred due to the cumulative misalignment of decision-making processes that resulted in single erroneous decisions. Therefore, early detection of drift allows for corrective action prior to contractual breaches; regulatory violations; or reputational damage occurring, reducing the financial liability associated with containing these failures and supporting stakeholder confidence.

Furthermore, drift detection enables organizations to delegate greater authority to agentic systems with increased confidence. Since managers are more likely to authorize decision-making authority to agentic systems when safeguards are provided to detect misalignment early, drift detection supports accelerated adoption of autonomy throughout procurement, logistics and inventory planning activities; thereby enabling organizations to realize performance improvements that may not have otherwise been achieved. In multi-agent supply chains, drift often manifests itself at the system level rather than individually at the agent level. While locally rational decisions made by multiple agents may collectively produce unstable global patterns of behavior (e.g. oscillatory inventory transfers, congestion amplification, or supplier switching cascades), drift detection must therefore consider aggregate behavior and coordination dynamics in addition to individual agent-level metrics (Zhang et al., 2021). It is essential for drift detection to include consideration of system-level phenomena in order to preclude emergent failure mechanisms.

Finally, drift detection provides organizations with additional support for regulatory defensibility by demonstrating ongoing continuous monitoring of autonomous behavior. Increasingly, regulators require documentation demonstrating that organizations proactively supervise AI-driven decisions rather than solely relying upon static compliance assertions. Therefore, the log files generated through drift detection provide organizations with this type of documentation through recording ongoing alignment checks and corrective actions. With regard to organizational learning, drift detection generates knowledge regarding how decision policy evolves in response to real-world conditions. Analysis of drift patterns provides insights to refine objectives, constraints, and reward structures. Through this feedback loop, the quality of governance and the performance of operational activities are continually strengthened over time. Ultimately, agent drift detection preserves the integrity of agentic supply chains by ensuring that autonomy is purposeful and not opportunistic. Through the integration of continuous alignment monitoring into system design, organizations can maintain adaptive intelligence without diminishing accountability; trust; or control.

Rollback and Safe Recovery Mechanisms

In order to provide an understanding of the role of rollback capabilities and safe recovery mechanisms in managing risk in agent-based supply chains, it is necessary to examine these two concepts together. As stated previously, autonomous execution of decisions must remain reversible; however, there is no guarantee that autonomous decision-making will always result in acceptable outcomes. Therefore, autonomous decision-making must include the ability to reverse decisions, or "roll back," when those decisions produce unacceptable results. The objective of this paper is to examine what rollback capabilities mean to agent-based supply chains and how they manage risk. Agent-based supply chains operate continuously, making interdependent decisions throughout procurement logistics, inventory positioning, and fulfillment. Because of this continuous nature of the agent-based supply chain, a short sequence of uncoordinated actions can spread rapidly through tightly-coupled networks. As a result, a key component of any strategy for containing risk in agent-based supply chains is to create rollback capability - a structural safeguard that recognizes that errors will occur but preserves the performance advantages of autonomous action.

However, rollback capability is not simply about reversing actions taken by the autonomous agent. Rather, it involves controlling the recovery of a supply chain to a stable condition following the occurrence of an error. For example, when inventory is redistributed or shipments are rerouted, the downstream feasibility conditions (e.g., capacity availability, contractual exposures, and service commitments) change. If a naive undo operation is used to reverse the actions taken by the autonomous agent, it could potentially reintroduce instability into the supply chain. Furthermore, many of the decisions made by autonomous agents are interdependent and sequential. Therefore, rollback mechanisms need to safely recover the system to a stable operating region rather than attempt to restore the system to its exact prior state. This approach reflects the fact that supply chains are path-dependent systems (Ivanov, 2017).

Rollback mechanisms operationalize the concept that delegated autonomy is still subject to the authority of the organization. Delegating autonomy to agents enables them to make decisions independently, but it does not relieve the manager of their ultimate responsibility for the consequences of those decisions. Rollback mechanisms provide a tangible way for managers to exercise that responsibility without completely eliminating autonomous decision-making. As such, rollback mechanisms allow organizations to implement a graded control structure where the organization can intervene surgically when risk levels are exceeded, while continuing to permit autonomous decision-making in other areas of the network.

As a result of the potential for autonomous routing and allocation decisions to rapidly build financial exposure, the business value of rollback capabilities is most evident in high-velocity logistics and fulfillment applications. Autonomous routing and allocation decisions can continue to increase financial exposure until some corrective action is taken, at which point the exposure will begin to decrease. Rollback can prevent compound interest from increasing the exposure beyond a certain point, thereby preventing the financial exposure from becoming systemic in nature. Rollback provides the financial protection needed to protect profit margins, service reliability and customer confidence. Thus, rollback represents a method to control financial risk, as well as to control technical risks.

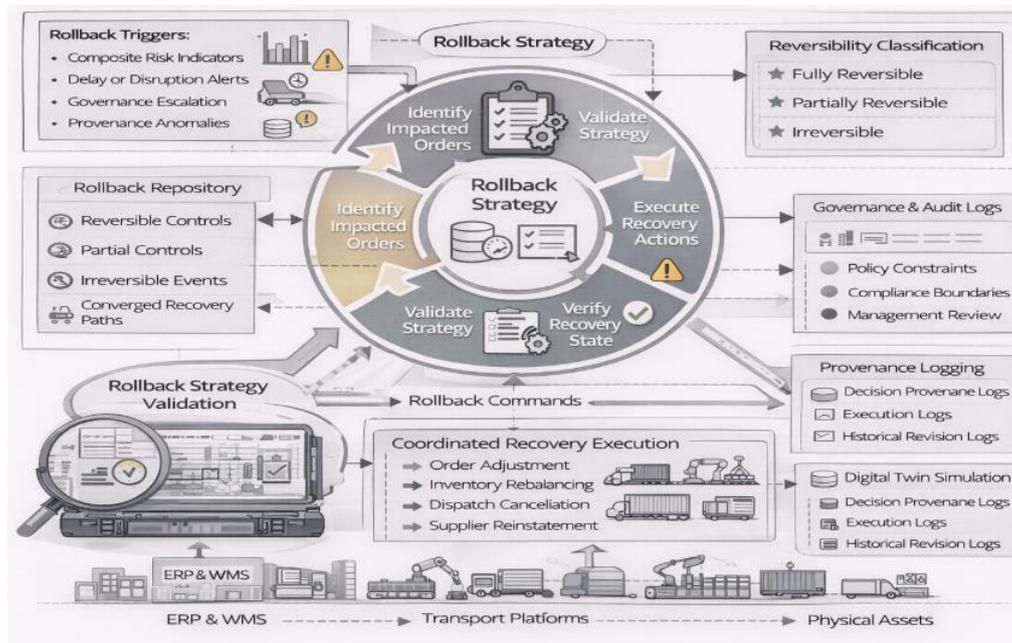
To effectively design rollback mechanisms, it is essential to classify decisions in terms of their reversibility, their impact, and their time-sensitivity. Some decisions, such as virtual inventory commitments or planning recommendations, are very reversible with little or no financial cost. Other decisions, such as physically shipping goods, canceling contracts, or terminating services, can become irreversibly committed once they are initiated. To ensure that high-risk decisions receive sufficient oversight, safe recovery mechanisms should incorporate provisions to delay the execution of high-risk decisions by longer periods of time, impose higher approval thresholds, or require additional review and validation. These provisions will help ensure that high-risk decisions do not impede the normal flow of routine operations. Rollback mechanisms are inherently integrated with provenance logging and digital twin technologies. Provenance logging technology is used to capture the specific sequence of decisions and state transitions that require rollback or containment. Digital twin technology is used to simulate alternative recovery scenarios to determine whether alternative recovery plans will successfully stabilize the supply chain or cause further disruption (Tao et al., 2019). Together, these technologies will ensure that recovery actions are grounded in the dynamics of the supply chain and not based upon uninformed intuition, which is particularly important in large and complex supply chains.

In addition to supporting the ability of organizations to maintain control over autonomous decision-making, rollback capability will also help to promote trust among internal stakeholders and external regulatory agencies that autonomous decision-making will not eliminate all controls over decision-making. When organizations have established clear recovery options, managers are more likely to delegate autonomous decision-making authority to agents. Similarly, when regulatory agencies see that autonomous systems are designed to be reversible, they are more likely to permit the use of autonomous decision-making systems. However, the safe recovery mechanism itself must operate within the confines of applicable laws, regulations, and contractual agreements. For example, the recovery actions required to correct the problem caused by an error may trigger a contractual obligation or penalty. Therefore, it is necessary to establish a framework for governance that includes recovery policies that are compliant with relevant laws and regulations.

In addition to providing an assurance mechanism for stakeholders that autonomous decision-making is reversible, rollback mechanisms can also be used to facilitate continuous improvement. Following each intervention, the organization can analyze what decision patterns, environmental conditions, and/or coordination failures led to the need for recovery. That knowledge can then be used to refine the organization's policy space, escalation thresholds, and governance constraints. Over time, the organization can develop greater resiliency by learning how to respond to errors more quickly and efficiently rather than trying to avoid errors altogether. When multiple agents are involved in a supply chain, the process of recovering from errors can become problematic. Specifically, when one agent reverses decisions that were made autonomously, it may render invalid the assumptions made by other agents. Therefore, safe recovery mechanisms need to coordinate the recovery efforts of all agents involved so that the recovery of one agent does not destabilize the entire supply chain. Coordinating recovery efforts among multiple agents can be done through coordinated rollback protocols that ensure consistency among agents who share resources, planning horizons, etc. (Li et al., 2020).

Finally, by embedding the principles of reversibility, stabilization, and coordinated intervention into the design of the architecture of autonomous supply chains, organizations can achieve high degrees of autonomy while at the same time protecting the interests of stakeholders, promoting resiliency, and establishing trust.

Figure 8: Rollback Strategy



Simulation Based Stress Testing

Simulation-based stress testing is a method for examining the reliability and safety of autonomous supply chain systems prior to their occurrence in actual supply chain operations. Unlike traditional supply chain testing methods that focus on average performance and historical data, simulation-based stress testing is a method for evaluating the behavior of autonomous decision-making systems during extreme low-probability high-impact events. During typical operation, autonomous decision-making systems can operate within acceptable bounds; however, under stress conditions, the systems may behave erratically or in unsafe ways. Simulation-based stress testing creates a controlled environment in which the behavior of autonomous decision-making systems can be examined during stressful events without risking physical damage to equipment or regulatory breaches. The controlled nature of simulation-based stress testing creates a "sandbox" in which autonomous agents can interact with simulated disruptions such as a port closure, a supplier bankruptcy, a transportation capacity reduction, or a sudden surge in demand.

Examination of autonomous decision-making systems under stress conditions will allow organizations to evaluate the stability of decisions, the effectiveness of coordination among agents, and compliance with governance requirements. Therefore, simulation-based stress testing transforms digital replicas of physical supply chain assets, referred to as "digital twins," into active tools for managing risks associated with autonomous decision-making systems. One of the key aspects of designing a stress testing program is developing valid stress scenarios. To be effective, stress scenarios need to be both extreme and realistic so that potential vulnerabilities in the supply chain are revealed, rather than being artificially induced. For example, in supply chain stress testing, scenarios should include correlated disruptions such as a supplier bankruptcy occurring simultaneously with increased regulatory oversight or a series of cascading logistics delays resulting from a major outage at a critical transportation facility. Additionally, the development of stress testing scenarios requires consideration of the interdependencies that exist between procurement, production, and delivery activities. Stress testing programs that use poorly constructed stress scenarios risk either exaggerating the vulnerability of the supply chain or providing false assurances about its capabilities. Therefore, the construction of rigorous stress testing scenarios is critical to achieving meaningful outcomes from stress testing programs (Shapiro et al., 2014).

Simulation-based stress testing also enables the evaluation of the coordination behavior of autonomous decision-making systems when subject to stress. Autonomous decision-making systems may exhibit emergent behaviors in multi-agent systems that do not occur in individual decision-making tests. When stressed, coordination mechanisms among agents may fail causing oscillatory inventory transfer patterns, congestion amplifications, and resource conflicts. Stress testing identifies these coordination failures, thereby providing an opportunity for

governance teams to redesign communication protocols and coordination mechanisms before deployment (Zhang et al., 2021). From a governance perspective, simulation-based stress testing provides an ex-ante validation of the boundaries of autonomy. Governance constraints and escalation thresholds may appear adequate under nominal conditions, but may be insufficient during stress conditions. Simulation testing determines if agents adhere to compliance constraints when faced with severe trade-offs, or if they improperly prioritize performance objectives. This insight enables governance parameters to be modified proactively, rather than reactively, after a failure. Thus, simulation-based stress testing operationalizes governance-by-design principles.

There are significant business benefits to simulation-based stress testing. Most notably, simulation-based stress testing reduces exposure to tail risk events that could threaten organizational viability. Supply chain failures often result from rare event combinations that were not previously considered. Stress testing enables organizations to systematically examine these event combinations and identify failure modes that may remain unknown. By identifying these failure modes and mitigating them in advance, organizations reduce their exposure to catastrophic supply chain disruptions and maintain stakeholder confidence (Ivanov, 2017). Additionally, simulation-based stress testing facilitates engagement with regulatory agencies and satisfies regulatory requirements for approval of autonomous systems. Regulatory bodies increasingly expect organizations that deploy autonomous systems to demonstrate that they have examined the behavior of those systems during adverse conditions. Providing regulatory agencies with simulation-based evidence of the assessment of autonomous systems' behavior during adverse conditions fosters greater regulatory trust and may be necessary in higher-risk industries, such as pharmaceuticals, defense, or critical infrastructure logistics. Thus, simulation-based stress testing serves as a compliance-enabling capability rather than solely a technical exercise.

To effectively conduct stress testing, organizations must consider the uncertainty and variability present in supply chain environments. Events rarely unfold in deterministic fashion and autonomous decision-making systems' behavior can vary between simulation runs due to random elements of the decision-making process. Therefore, effective stress testing requires multiple simulation runs and statistical analysis of the results to determine the robustness of the system rather than analyzing individual run trajectories. This probabilistic approach to assessing supply chain risk is consistent with the experience of supply chain managers and helps to avoid overconfidence in isolated simulation results (Shapiro et al., 2014). Simulation-based stress testing enables organizations to comparatively evaluate alternative decision policies and governance structures. Organizations can simulate how different constraint values, escalation thresholds, and coordination architectures behave under the same stress scenarios. The controlled comparison enabled by simulation-based stress testing supports evidence-based design of governance by illustrating which configurations achieve the optimal balance between performance, resilience, and compliance. Due to the operational risks associated with live experimentation, obtaining this type of evidence-based information is difficult using other methodologies. A simple mathematical representation of stress testing can be formulated using the concept of expected loss given the stress scenarios. If S denotes a set of stress scenarios and L denotes a loss function representing the costs associated with service degradation or compliance breaches resulting from policy π :

$$\mathbb{E}[L | \pi] = \sum_{s \in S} P(s) \cdot L(s, \pi)$$

This formula illustrates the principle that the evaluation of governance and decision-making policies under stress is dependent upon both the probability of each stress scenario occurring and the impact of each scenario. Since simulation provides estimates of these two variables, simulation-based stress testing enables informed decision-making about governance policies (Shapiro et al., 2014). Simulation-based stress testing also facilitates organizational learning and preparedness for potential crises. Governance teams and decision-makers gain an understanding of the complex failure dynamics that cannot be intuited from the experience alone. This understanding enables better preparation for responding to crises and a deeper understanding of how systems behave under duress. As the number of simulations increases, stress testing becomes an integral part of organizational culture and enhances resilience. Simulation-based stress testing transforms risk management in agentic supply chains from a reactive response to a proactive validation of the reliability and safety of autonomous decision-making systems. Rather than relying on post-failure corrective action to address issues that arise in supply chains, simulation-based stress testing allows organizations to proactively validate the reliability

and safety of autonomous decision-making systems, and therefore, to incorporate resilience into system design (Ivanov, 2020).

Resilience Under Adversarial Conditions

Developing mechanisms to enable agentic supply chains to operate under adversarial conditions requires the development of architectures and governance structures to facilitate resilient operation (Goodfellow et al., 2015). Supply chains are becoming increasingly autonomous and data-driven, therefore, they will become increasingly attractive to individuals seeking to take advantage of the systems for financial, competitive advantage, or political leverage (Goodfellow et al., 2015). Autonomous systems can react to manipulated inputs mechanically and therefore potentially exacerbate the negative effects of adversarial attacks (Goodfellow et al., 2015). An explicit architecture and governance structure are needed to enable autonomous systems to operate safely in environments hostile to their interests (Goodfellow et al., 2015).

Adversarial conditions can exist in several forms within supply chain operations, including; Data poisoning (the act of contaminating data), Falsification of Demand Signals, Spoofing Shipment Telemetry, Compromised Supplier Information, and Interference with Logistics Operations (He & Zhang, 2023). In contrast to stochastic disruptions that occur randomly, adversarial disruptions tend to be both strategic and adaptable (He & Zhang, 2023). Adversaries adapt their tactics and strategy based upon the response of the affected organization and therefore, to effectively design and implement agentic supply chains, expect to receive deceptive or malicious inputs.

A fundamental aspect of developing agentic supply chains capable of operating under adverse conditions is the assessment of the credibility of information sources in real-time. Supply chain agents receive data from a variety of sources, including internal enterprise systems, external suppliers, and partner systems. However, in an adversarial environment, it is reasonable to anticipate that some portion of the input received will be corrupted. To mitigate this issue, resilient agentic systems utilize redundancy, cross-validation, and consistency checking to verify if the received data is consistent with historical trends and other relevant data (Chandola et al., 2009). For instance, if an increase in demand occurs quickly with no commensurate increase in downstream consumption, it is most likely that the increase in demand was artificially created. If the system identifies that the received input is most likely false, the system can choose to avoid reacting aggressively to the input.

From a governance perspective, the need to develop and implement mechanisms to provide for resiliency in the presence of an adversary will necessitate that agents react conservatively when they do not believe the inputs to the decision-making system (Goodfellow et al., 2015). When the system determines that the inputs to the decision-making system are unlikely to be trusted, the agents must revert to established policies that place emphasis on safety and compliance rather than aggressive optimization and expansion of operations (Goodfellow et al., 2015). Examples of how this can manifest itself include limiting the degree of change made to operational parameters, limiting the frequency of escalations, or reverting back to established baseline policies. Establishing and implementing governance structures that define fallback behaviors will assist in ensuring that the system responds in a predictable manner to an adversary's attempts to disrupt its operations, as opposed to responding in an ad-hoc manner (Yuan et al., 2019).

The implications for businesses and organizations in developing mechanisms for resiliency in the presence of an adversary are substantial, primarily because the disruption of supply chains can result in severe financial and reputational consequences for organizations (Goodfellow et al., 2015). Misallocating inventory, or diverting shipments due to an adversarial actor can lead to widespread failures of service, customer dissatisfaction, and potentially regulatory action. By developing mechanisms for resiliency in their systems, organizations can preserve their revenue streams and protect their brand's reputation and trust, regardless of the environment that is hostile to their interests (Goodfellow et al., 2015). Resiliency has progressed from being a simple organizational defense mechanism to a competitive necessity.

The concept of resiliency in the context of adversarial conditions has additional implications relative to the dynamics of coordinating multiple-agents in a supply chain. Coordinated attacks can exploit the interactions between multiple-agents to create instability, such as synchronized inventory oscillations or congestion cascades,

by inducing agents to collectively amplify their own instability (Mirsky et al., 2018). Resilient systems must monitor not only the behavior of each individual agent, but also the emergent patterns of behavior in the system that may represent a coordinated effort to manipulate the system. Identifying such patterns enables organizations to confine the impact of the attack at the system-level as opposed to identifying and responding to each individual agent separately (Mirsky et al., 2018).

Simulation-based adversarial testing provides another means to evaluate the effectiveness of resilient systems and to simulate an adversary adapting their tactics based on the system's responses to prior attempts to manipulate the system. Through simulation-based testing, organizations can identify vulnerabilities in the systems utilized to validate signals, coordinate activity among agents, and establish thresholds for governance and compliance. Simulation-based testing offers organizations the opportunity to proactively make informed decisions regarding defensive design options that may be challenging to identify utilizing benign testing methods alone (Carlini & Wagner, 2017).

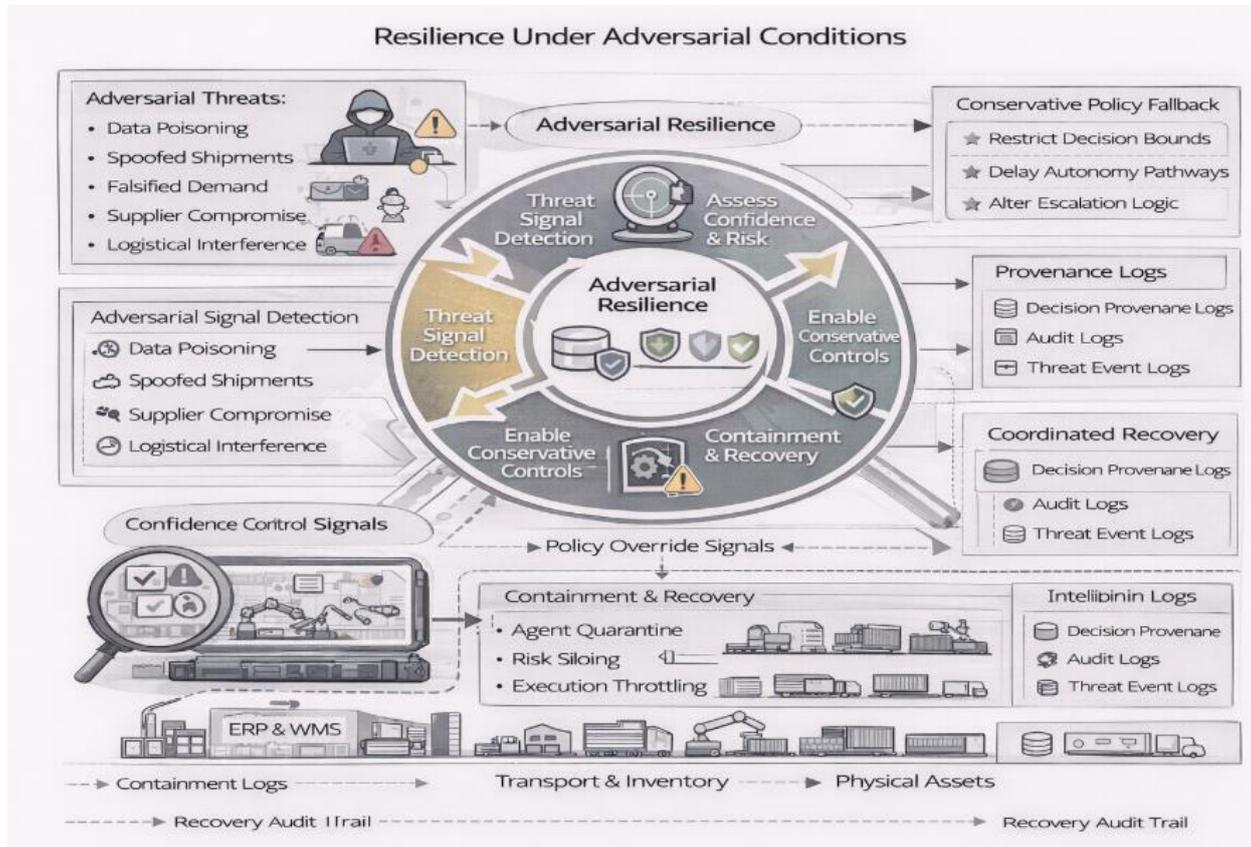
In the context of resiliency in adversarial conditions, the capability to rapidly contain and recover from disruptions is also involved. Although robust detection mechanisms can detect the majority of disruptions caused by an adversary, there will be instances where the system cannot entirely prevent an adversary from causing damage. Rapidly containing the scope of decision-making for agents impacted by an adversary, and rapidly restoring normal operation to the system, is vital for reducing the total amount of damage caused by an adversary. Organizations must pre-define and audit their recovery processes to ensure that the actions taken during the recovery process are compliant with organizational policies and regulatory requirements (Kshetri, 2018).

Finally, from a regulatory and societal perspective, the development of mechanisms for resiliency in the context of adversarial conditions supports the objective of protecting critical infrastructure. Critical infrastructure includes a broad array of essential services including food production and distribution, healthcare delivery, and energy production and delivery. Regulators are now requiring organizations to demonstrate their ability to resist disruption through the utilization of autonomous systems, in addition to demonstrating compliance with regulatory requirements. Demonstrating an organization's ability to resist disruption through the utilization of autonomous systems can positively impact the regulator's assessment of the organization's suitability to operate in a domain in which autonomous systems are employed, and may also positively influence public trust in the organization (Goodfellow et al., 2015).

Data sovereignty and decentralized architectures can also provide benefits to organizations attempting to build resiliency into their systems. Decentralizing decision-making authority and limiting coordination across geographic areas can limit the extent to which an adversary's manipulation of input signals in one region of the system can affect decision-making globally. Limiting coordination while preserving regional autonomy can reduce the systemic risk associated with the employment of autonomous systems. Thus, the establishment of resiliency in systems through the use of decentralized architectures aligns with the objectives of data sovereignty (Tao et al., 2019).

Ultimately, the development and implementation of the mechanisms required to ensure that agentic supply chains continue to operate in an environment that is antagonistic to their interests is critical to achieving the objective of retaining the trust and reliance of stakeholders on the systems developed. By including mechanisms for resiliency in systems through detection, conservative decision-making, containment, and recovery, organizations can develop autonomous systems that degrade gradually in the case of an attack, rather than failing catastrophically (Madry et al., 2018).

Figure 9: Resilience Under Adversarial Conditions



Human Oversight and Organizational Integration

Human oversight is a necessary component to support both the legitimacy and sustainability of autonomous agent-based supply chains. With agent-based systems taking over operational control for tasks including purchasing, logistics, inventory planning and fulfillment; human roles transition from being directly involved in the decision-making process to serving as supervisors in a governance position. While human oversight continues to be an essential part of the decision-making process in an agent-based supply chain environment; the type of oversight needed changes dramatically. Instead of being transaction-oriented, oversight needs to be interpretive-strategic-corrective oriented. As a result, the success of agent-based supply chains will depend upon the technical architecture of the agents as well as how the human oversight is institutionalized within an organization's structure and processes.

There exists a significant difference in the way oversight is implemented within agent-based supply chains depending on whether it is based on either a human-in-the-loop or a human-on-the-loop model. A human-in-the-loop model is based upon the need for direct human approval or intervention before the autonomous decision is made by the agent. This model offers a high degree of control but introduces latency and severely limits scalability. As a result, this model is not suitable for high-frequency supply chain decision making. A human-on-the-loop model allows for humans to serve as supervisors who monitor agent behavior; can choose to intervene selectively and/or modify governance parameters as opposed to approving each and every action made by the agent. Given the need to achieve scalable supply chain operations while still maintaining accountability; the human-on-the-loop model is essential in the context of agent-based supply chains. The choice between these two models is primarily dependent on an organization's risk tolerance; regulatory requirements and level of operational criticality.

Effective implementation of human-on-the-loop oversight relies heavily on the availability of sophisticated monitoring and interpretive interfaces that allow for supervisors to comprehend autonomous agent behavior without the need for micro-management (Kaber, 2018). Managers responsible for supply chain activities must be capable of observing patterns, trends and anomalies that exist across multiple decision streams versus viewing

each action in isolation. To accomplish this objective, supply chain managers require access to dashboards and analytical tools that surface governance relevant signals including; drift indicators, compliance stress levels, escalation frequencies, etc. If such visibility is not available; then human oversight becomes symbolic as opposed to providing substantive oversight capabilities. Therefore, effective oversight is contingent upon having technological systems that provide supervisors with actionable managerial insight regarding autonomous agent behavior.

Accountability, as it relates to decision making, remains paramount even though decision making responsibilities have been delegated to autonomous agents. Organizations cannot assign failures to algorithms without destroying governance credibility (Norman, 1990). Accountability must be redefined so that managers are held accountable for the design and configuration of autonomous systems and their overall governance as opposed to holding individuals accountable for specific decisions. Such a redefinition of accountability is consistent with current principles of corporate governance where executives are accountable for systems of control as opposed to all operational acts. For example, in supply chain operations; accountability rests with those responsible for defining policy constraints and escalation rules and ensuring the overall effectiveness of oversight.

The redistribution of accountability has significant legal and regulatory implications. Courts and regulators continue to seek identifiable human accountability even when decisions are being made using algorithms (Shneiderman, 2020). The ability to clearly define managerial accountability for autonomous systems provides this identifiable anchor. Therefore, organizations must formally define roles such as autonomous system owner, governance steward and escalation authority. These defined roles assist in ensuring that responsibility for autonomous agent behavior is explicit and not diffuse. Clarity is essential for obtaining regulatory acceptance and internal governance discipline.

Another key aspect of integrating humans into agent-based supply chain operations is establishing trust calibration. Excessive trust in autonomous systems may lead to complacency and delayed intervention, whereas insufficient trust may result in excessive override of autonomous decisions thereby negating the benefits of autonomy (Hoff & Bashir, 2015). Trust calibration refers to developing human confidence in actual system capability and reliability. In supply chain operations, trust calibration is developed through transparency, explainability, and consistent system performance. When autonomous agents perform in a predictable manner and provide logical explanations for their decisions, managers develop calibrated trust that supports effective oversight.

Like many aspects of integrating humans into agent-based supply chain operations, trust calibration is dynamic rather than static. As autonomous agents learn and adapt to changing environments, human trust must be continually recalibrated (Dzindolet et al., 2003). Even if an autonomous agent behaves in a manner that is technically correct, unexplained unexpected behavior may cause human trust to deteriorate. On the other hand, continued reliable performance during times of stress increases human confidence. Therefore, organizations must invest in communication methods that inform human supervisors about system learning and adaptation, system limitations, and potential risks. An ongoing dialogue between humans and machines is essential for successful, stable integration.

Finally, successful embedding of agent-based supply chains is significantly dependent on the ability of an organization to effectively manage the necessary changes to job roles, decision authority, and performance metrics. Without intentional organizational change management, autonomous agents may create resistance, anxiety, and/or role confusion among supply chain personnel (Merritt et al., 2013). Successful change management redefines autonomy as an enhancement to human expertise in oversight and governance. Training programs must provide managers with the skills required to oversee autonomous agents instead of executing routine decisions.

In addition to providing training, change management must also involve aligning incentives and performance metrics with the new roles created in agent-based supply chain environments. Traditional performance metrics used to measure the efficiency of transactions may no longer accurately measure the contribution of managers in agent-based environments (Madhavan & Wiegmann, 2007). Organizations must revise their evaluation

frameworks to include performance metrics that recognize effective oversight, risk management, and governance quality. Aligning manager incentives and performance metrics ensures that managers are encouraged to interact positively with autonomous agents rather than resisting them. In supply chain environments where autonomous agents reshape workflows, aligning performance metrics and incentives is critical to creating an organizational culture that accepts and supports autonomy.

In addition to supporting organizational change and performance improvement, human oversight is also necessary to support the ethical and societal obligations inherent in agent-based supply chains. Decisions related to supplier selection, labor practices, environmental issues, and service priorities carry normative implications that cannot be completely delegated to algorithms (Hancock et al., 2011). Governance bodies comprised of humans must establish ethical guidelines and ensure that autonomous agent behavior is aligned with organizational values and societal expectations. Therefore, oversight extends beyond the management of operational risks to include reputational and ethical considerations that impact the long-term viability of businesses.

Ultimately, the integration of human oversight will determine whether agent-based supply chains are viewed as trustworthy organizational systems or isolated technical artifacts (Shneiderman, 2020). Autonomous agents without oversight may become opaque and lose credibility, while oversight without autonomy reduces scalability and response time. Through the use of human-on-the-loop supervision, clear definitions of managerial accountability, calibrated trust, and structured change management, organizations can reconcile human judgment with machine execution. Ultimately, this reconciliation will enable agent-based supply chains to realize performance improvements while maintaining the attributes of governance, accountability, and social acceptability.

Enterprise Integration and Deployment Considerations

The degree to which enterprise integration and deployment issues impact whether agentic supply chain architectures can evolve into operational systems embedded in organizations (Xu et al., 2018) will depend upon how well the enterprise integration platform can adapt to the evolving needs of the organization. While agentic intelligence has the promise to produce adaptive decision-making and autonomous actions, the practical utility of this intelligence is dependent upon seamless integration with the existing enterprise integration platform, including all transaction, inventory flow, and operational control elements. Most organizations have developed and are using complex systems composed of Enterprise Resource Planning (ERP), Supply Chain Management (SCM) and Warehouse Management Systems (WMS) software, which contain the years of process logic, compliance rules, and financial controls. Therefore, the agentic systems must seamlessly integrate with these systems, without disrupting core operations or violating established governance structures.

The most significant integration requirement is integration with ERP systems because ERP systems serve as the authoritative system of record for financial transactions, procurement contracts, and master data (Benlian et al., 2009). The agentic supply chain must be able to interact with the ERP systems to perform purchasing decisions, manage supplier relationships, and capture the accurate costs associated with those decisions. This integration cannot be superficially done; autonomous decisions made without going through financial controls will undermine auditability and fiscal accountability. Therefore, agentic execution must be mediated by ERP interfaces, which enforce posting rules, approval hierarchies, and reconciliation logic. This will ensure that autonomy will always operate within the financial governance structure of the organization. Supply Chain Management (SCM) platforms are another key layer of integration required because SCM platforms manage and coordinate planning, forecasting and execution activities across procurement, production, and distribution (Wamba et al., 2015). The agentic systems typically augment or replace traditional planning logic within these platforms by adding adaptive policies that react to real-time events. Therefore, the integration must support bi-directional interaction between the agentic system and SCM platforms, so that agentic system decisions update plans and execution status while receive constraints, forecasts and performance feedback. If this bi-directional interaction does not occur, then there will be discontinuity between strategic planning and autonomous execution instead of having parallel control structures that compete for authority.

There are unique integration requirements for Warehouse Management Systems (WMS) because WMS systems operate at a high-frequency, low-latency within physical environments (Lee & Lee, 2015). Autonomous decisions regarding picking, allocation, slotting, and replenishment must be synchronized with warehouse execution logic to prevent operational conflicts. Additionally, agentic systems must respect deterministic constraints of warehouse operations such as equipment availability, labor scheduling, and safety rules. Thus, the integration of agentic systems with WMS systems must be carefully orchestrated to ensure that autonomous decisions are translated into executable tasks that align with physical workflows.

Application Programming Interfaces (API)-mediated orchestration represents the predominant mechanism for integrating agentic systems with enterprise platforms (Jamshidi et al., 2018). Rather than relying on direct database access or hard-coded logic, agentic execution should utilize standardized APIs to encapsulate business rules and validation logic. Utilizing API-mediated orchestration preserves system integrity and enables incremental deployment without extensive re-engineering. Furthermore, API-mediated orchestration facilitates modularity, allowing organizations to introduce agentic capabilities in a gradual manner across different supply chain functions.

From a governance standpoint, utilizing interface-based integration increases auditability and control (Panetto et al., 2019). Every autonomous action taken by an agentic system passes through enterprise interfaces, which can log, validate and constrain the action according to policy. As such, the mediation of agentic execution via enterprise interfaces ensures that agentic execution remains observable and compliant even as decision logic evolves. Additionally, utilizing interface-based integration enables organizations to selectively throttle or suspend autonomy by controlling interface permissions rather than changing the core algorithms used in decision logic. Throttling or suspending autonomy through interface permissions is necessary to control risk when deploying agentic systems early in their lifecycle.

Hybrid cloud and edge execution architectures play a crucial role in achieving a balance between scalability, responsiveness and compliance in agentic supply chains (Tao et al., 2019). Cloud environments provide elastic computing resources for policy learning, simulation and global coordination. Edge environments located near operational assets provide low-latency decision-execution, where timeliness is critical, such as warehouse control or transportation routing. Hybrid architectures allow agentic systems to distribute intelligence across layers while adhering to data locality and latency constraints. Distribution of intelligence across layers is necessary to achieve real-time supply chain responsiveness.

Latency is an especially significant consideration for agentic execution since latencies between decision-generation and decision-execution can degrade performance or cause instability. In cases where timeliness is critical, such as transportation routing or warehouse control, decisions must be executed within strict temporal bounds (Ghofrani et al., 2018). Hybrid execution architectures reduce latency by placing decision-logic close to the point-of-action while maintaining coordination with higher-level intelligence. However, achieving this balance requires careful design of the overall system architecture, rather than simply deploying the architecture. Another deployment challenge for agentic supply chains is scalability since agentic supply chains must operate across large networks with potentially thousands of decision-points (Panetto et al., 2019). Integration architectures must support horizontal scalability without creating bottlenecks at the enterprise interfaces. To accomplish this, integration architectures must employ asynchronous communication patterns, event-driven processing, and robust interface designs. Achieving scalability is a performance concern, but it is also a governance concern, since overloaded systems may bypass controls or compromise auditability under load. A simple expression for the latency-constrained execution problem can be represented as follows:

$$t_c + t_i \leq T$$

where T is the maximum allowable latency, t_c is the time to generate the decision, and t_i is the time to execute the decision.

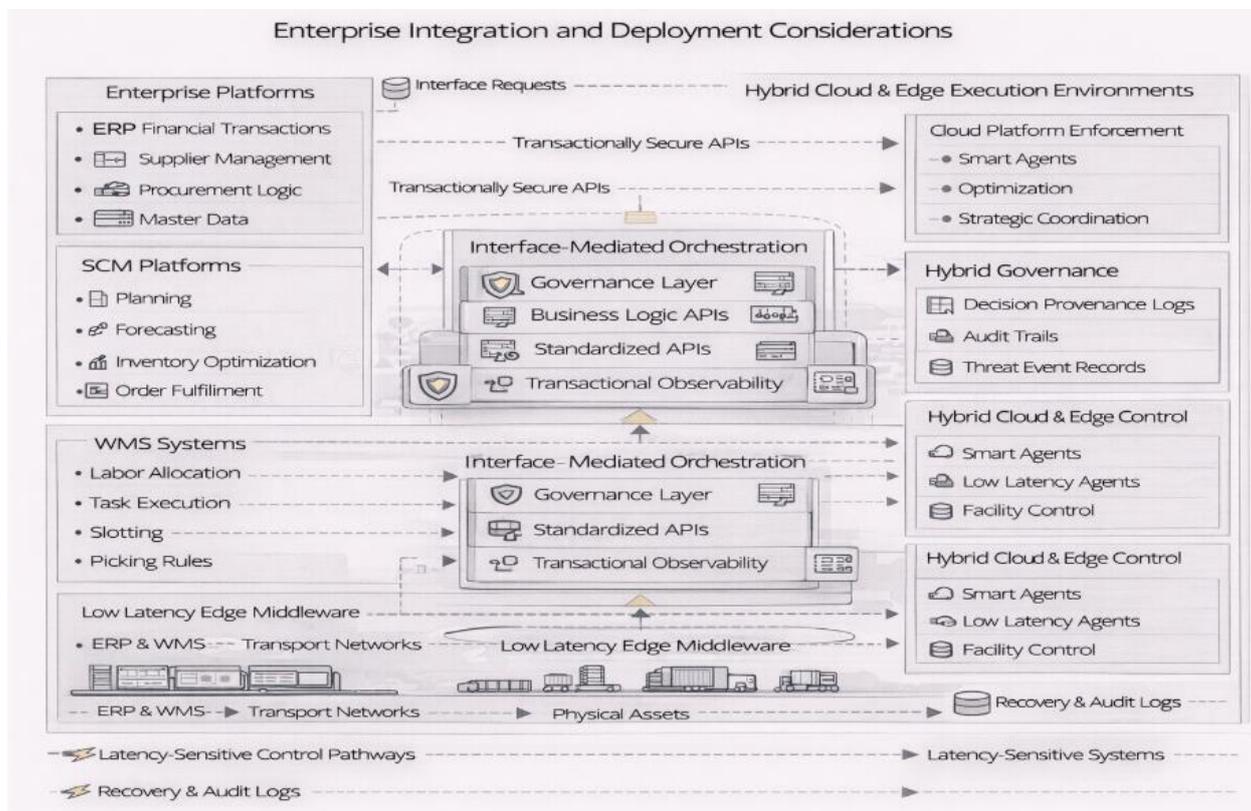
This formulation illustrates that the integration architecture must take into account both the time to compute the decision (intelligence) and the time to mediate the decision through the enterprise systems (execution) to meet

the operational deadlines. In many supply chain operations, failing to meet the constraint illustrated above may result in obsolete or disruptive decisions (Ivanov & Dolgui, 2020).

In addition to affecting the latency of the decision execution, enterprise integration affects an organization's preparedness to deploy agentic supply chain capabilities. An organization should introduce agentic capabilities incrementally, starting with low-risk decision-domains and increasing over time as confidence grows (Hosseini et al., 2019). The integration architectures that facilitate modular deployment enable this incremental introduction of agentic capabilities without causing instability in the core systems. This incremental introduction of agentic capabilities is consistent with introducing agentic capabilities in conjunction with organizational change management and risk tolerance.

Ultimately, the deployment of agentic supply chains will determine whether agentic supply chains produce long-term benefits or remain isolated demonstrations. Successful deployment of agentic supply chains is contingent upon integrating the agentic capability with the existing enterprise platforms, thereby embedding the autonomy in the existing control structures rather than in separate, parallel control structures. Organizations that successfully design robust interfaces, hybrid execution architectures, and scalable integration architectures will be able to bridge the conceptual agentic architectures with the realities of operational environments and realize the potential of agentic supply chains (Ivanov & Dolgui, 2020).

Figure 10: Enterprise integration considerations



Evaluation Metrics and Benchmarking Frameworks

Metrics for evaluating the performance of an autonomous supply chain should be viewed as tools of governance, and therefore, not simply as performance evaluation scorecards (Lipton, 2018) since traditional supply chain metrics are primarily focused on accuracy, cost, and service levels. For example, metrics for an autonomous supply chain would need to evaluate whether the decision authority that has been delegated to an agent is being used reliably, transparently, compliantly and resiliently over time. Since autonomous supply chains will continue to adapt continuously at scale, metrics for evaluating their performance will need to measure both behavioral stability and structural risk accumulation over time, as well as the level of institutional trustworthiness. Therefore, metrics for evaluating autonomous supply chains do not just determine how the system is evaluated,

but also how the system evolves, as the learning dynamics of the system respond to what is measured and rewarded.

Decision Consistency Index: The decision consistency index measures the degree to which an agent produces similar actions when provided with substantially similar supply chain states, taking into account factors such as demand patterns, capacity constraints, regulatory context and contractual obligations (Busoniu et al., 2008). For example, in a supply chain setting, this metric would assess whether procurement, sourcing, allocation, routing, or inventory position decisions remain consistent when the underlying conditions are essentially the same. If there is excessive variability in the decisions produced when compared to the number of materially equivalent states, it could suggest that the decision-making policies are brittle or overly reactive and potentially exacerbate operational noise into systemic volatility. Consistent decision-making does not equate to rigid repetition; rather, it means that agents produce decisions that vary only when material differences in states have occurred. By tracking the consistency of decision-making over time, organizations can evaluate whether the autonomous decision-making process is reliable and trustworthy enough to facilitate trust and coordination across connected supply chain functions.

Policy Entropy Level: The Policy Entropy Level measures the variance in the action choice set that an agent considers or chooses under a given supply chain state. In operational terms, this metric reflects whether an agent oscillates excessively between alternatives, or whether an agent converges too quickly on limited decision patterns (Busoniu et al., 2008). Excessive entropy in supply chain decision-making processes can lead to indecision, increased frequency of changing suppliers, unstable routing selections, and inconsistencies in allocating inventory. All of these scenarios create additional friction and coordination costs in executing supply chain activities. On the other hand, extreme low entropy can indicate that an agent has lost its ability to adapt to regime shifts and unforeseen disruptions due to excessive confidence in its current decision-making processes. The longitudinal tracking of entropy can allow organizations to monitor whether the learning dynamics in the autonomous decision-making process are becoming stabilized properly, or whether they are trending towards pathological extremes.

Behavioral Drift Magnitude: Behavioral Drift Magnitude represents the total deviation of an agent's decisions from the baseline definitions established by governance over time (Zhou & Li, 2020). In the context of supply chains, Behavioral Drift Magnitude can detect slow-moving changes such as an increasing willingness to accept concentrated risk by suppliers, gradual erosion of compliance margins, and systematic biases in favor of cost minimization at the expense of service reliability. Drifts are particularly damaging because they frequently go unnoticed in short-term performance metrics, yet accumulate structural risks over time. By quantifying the magnitude of drift, organizations can identify potential issues and implement corrective actions before they become operationally apparent and compromise the performance and governance of adaptive autonomous systems.

Coordination Stability Score: The Coordination Stability Score assesses whether a group of agents maintains coordinated, non-oscillatory coordination of shared supply chain resources, such as inventory pools, transportation lanes, production capacity, and supplier commitments. In practice, poor coordination stability results in repetitive inventory shuttle movements, congestion cascades, or simultaneous switches in suppliers, which negatively impact global performance. This metric assesses whether individual, local autonomy contributes to system-wide coherence or to emergent instability. High coordination stability implies that decentralized decision-making respects common constraints and produces predictable collective behavior necessary for the scaling of agentic supply chains.

Unexpected Override Frequency: Unexpected Override Frequency assesses how often human supervisors intervene outside of predefined escalation pathways to either correct or stop autonomous decision-making. In supply chain operations, high frequencies of unexpected overrides typically signify either that the autonomous agent(s) are producing unreliable decisions, or that the governance constraints on the agent(s) are poorly calibrated. This metric is a proxy for trust erosion, as operators tend to intervene when autonomy is unpredictable or opaque. A sustained low override frequency indicates that autonomous agents are behaving predictably within

expected bounds, and that the governance design successfully anticipates edge cases and enables humans to remain in a supervisory role versus a reactive role.

Decision Provenance Coverage Rate: Decision Provenance Coverage Rate represents the percentage of autonomous supply chain decisions for which complete provenance records exist, including the source data inputs, evaluated constraints, alternative actions considered, and the actual execution outcomes (Ribeiro et al., 2016). In regulated supply chains, any gaps in provenance undermine audit-defensibility, regardless of the quality of the outcomes. High provenance coverage ensures that decisions remain reconstructible and institutionally accountable, facilitating regulatory reviews, contractual disputes, and internal governance oversight. As such, this metric is fundamental to treating autonomy as an auditable organizational process rather than as an opaque algorithmic activity.

Context Fidelity Score: Context Fidelity Score assesses whether the transient operational conditions (such as supplier availability, transportation capacity, regulatory status, and demand volatility) are accurately captured at the time of the decision. Supply chain decisions are heavily context-dependent, and inaccurate or incomplete context renders audit records misleading. High fidelity ensures that post-hoc analysis accurately reflects the operational conditions under which decisions were made, rather than being based on assumed conditions. This metric is crucial for establishing meaningful accountability and learning in agentic environments.

Governance Constraint Logging Completeness: Governance Constraint Logging Completeness assesses whether all applicable governance rules were evaluated and logged for every autonomous decision. In supply chains, this includes trade compliance checks, contractual service commitments, emissions thresholds, risk limits, and data sovereignty rules. Completeness signifies that governance was enforced systematically, and not implicitly or selectively. This metric provides evidence that the autonomy mechanism internalizes governance logic, and is not reliant on external enforcement mechanisms.

Temporal Dependency Preservation Index: Temporal Dependency Preservation Index measures whether audit logs preserve the correct order and causality among sequences of decisions. In many supply chains, outcomes arise from interacting decisions over time, rather than discrete actions. The preservation of temporal dependencies enables root-cause analysis and defensible causal attribution. Without this metric, audit trails devolve into disconnected events unable to explain the systemic behavior of the supply chain.

Replayability Availability Ratio: Replayability Availability Ratio measures the proportion of decisions that can be recreated and replayed within a digital twin environment. Replayability enables experiential audit validation through enabling reviewers to view decision behavior dynamically. A high replayability availability ratio enables deeper learning governance refinements and enhanced regulatory confidence, converting audits from static inspections into behavioral verifications.

Violation Rate per Decision Class: Violation Rate per Decision Class measures the frequency of regulatory or contractual violations per decision type (e.g., sourcing, routing, allocation). Segmentation is necessary, as aggregate compliance rates can mask localized risk. This metric enables targeted governance intervention and highlights which decision classes are most in need of tighter constraints, or better escalation logic.

Near Miss Frequency: Near Miss Frequency captures how often autonomous decisions approach compliance boundaries without crossing them. In supply chains, frequent near-misses indicate that governance is stressed, and that the organization is increasingly exposed to risk, even when no formal violations have occurred. Monitoring this metric enables proactive policy tightening prior to failures occurring, thus maintaining compliance margins in dynamic environments.

Escalation Appropriateness Ratio: Escalation Appropriateness Ratio assesses whether decisions are escalated when appropriate and not escalated inappropriately. In supply chains, inappropriate escalation can either delay execution, or expose the organization to risk. This metric assesses the quality of calibration of governance thresholds and determines whether human oversight enhances autonomy.

Jurisdictional Compliance Stability: Jurisdictional Compliance Stability measures the consistency of compliance performance across regions with varying regulatory regimes. Global supply chains must avoid uneven risk distribution, where certain jurisdictions incur disproportionately higher violation rates. Stability indicates that sovereign-aware governance has been embedded within the agentic decision-making process.

Override Justification Completeness: Override Justification Completeness assesses whether human overrides of autonomous decisions are accompanied by rationales that align with governance rules. This metric ensures that human interventions remain auditable, and not arbitrary, thereby ensuring institutional accountability, even when autonomy is curtailed.

Cost Volatility Reduction Index: Cost Volatility Reduction Index measures the decrease in cost variability, rather than simply the mean cost. In supply chains, volatility drives risk, operational stress, and managerially imposed burdens. Autonomous systems that reduce cost profile volatility provide superior business value by improving predictability and enhancing financial planning capabilities, even if the mean costs remain unchanged.

Service Level Stability Metric: Service Level Stability Metric assesses the consistency of service performance over time, rather than solely focusing on peak performance achievement. Stable service builds customer trust and reduces firefighting costs. Agentic systems must demonstrate resilience across cycles to justify autonomy.

Inventory Turnover Efficiency: Inventory Turnover Efficiency measures how effectively agents balance product availability with working capital utilization. Excessive turnover increases disruption risk, while inadequate turnover results in tying up working capital. This metric reflects disciplined operational control, rather than aggressive optimization.

Decision Cycle Time Compression: Decision Cycle Time Compression measures the decrease in time between signal detection and execution. In supply chains, faster cycles enhance responsiveness; however, compression must not negatively affect governance quality. This metric is meaningful only when evaluated together with compliance and audit metrics.

Managerial Load Reduction Index: Managerial Load Reduction Index measures the reduction in manual decision workload. Effective autonomy shifts human effort from execution to oversight, enabling strategic focus and organizational scalability.

Time to Anomaly Detection: Time to Anomaly Detection measures how rapidly deviations from expected patterns are detected. Rapid detection of anomalies prevents cascading disruptions, and reflects an effective monitoring infrastructure.

Time to Containment: Time to Containment measures how quickly autonomous behavior is contained after risk detection. Faster containment limits the blast radius and maintains operational continuity.

Recovery Duration: Recovery Duration measures the time required to restore stable operations after a disruption. A shorter recovery duration indicates effective rollback and recovery integration.

Performance Degradation Under Stress: Performance Degradation Under Stress measures how significantly the cost, service, or compliance deteriorate when subjected to adverse operating conditions. Resilience in autonomous decision-making processes is reflected in graceful degradation.

Structural Dependency Concentration: Structural Dependency Concentration assesses reliance concentrations of suppliers, routes, or regions created by agent decisions. High concentrations create hidden fragilities, even when performance appears optimal. Monitoring this metric prevents brittle optimization and ensures long-term resilience.

Frameworks for Benchmarking: Benchmarking frameworks integrate these metrics across time, configuration, and organizational unit boundaries (Ribeiro et al., 2016). Baseline benchmarking compares agentic systems to historical rule-based approaches. Configuration benchmarking assesses the comparative effectiveness of

different governance designs. Cross-regional benchmarking identifies disparate risk exposures (Zhou & Li, 2020). Composite trust benchmarking assesses balanced performance across all of the metric families.

$$T = w_r R + w_a A + w_c C + w_b B + w_s S$$

where trust emerges from balanced governance rather than isolated efficiency.

Strategic and Business Implications

Implications of agentic supply chains arise from the shift from episodic planning and human-mediated coordination to continuous governance of autonomous operations at machine speed (Waller & Fawcett, 2013). When autonomous decision-making is designed with governance mechanisms that include explicit constraint of decision-making; auditability; compliance; alignment with regulations; and resilience safeguarding; the supply chain becomes more than an operational cost center and develops into a strategic capability that defines market responsiveness, customer experience and risk posture. The reframed positioning of agentic supply chains as a durable source of competitive value stems from changing how organizations perceive their environment, identify and execute against threats, disruptions, and opportunities (Teece, 2007). Competitive advantage does not come from autonomy but from autonomous decision-making that is governable and trustworthy among all stakeholders including internal stakeholders; external partners; and regulators.

Governing autonomy provides a competitive advantage through compressed decision cycles while maintaining control of operations; thereby providing organizations with the capability to respond to changes in demand; disruptions; and market opportunities before competitors using manual decision loops (Kache & Seuring, 2017). Competitive differentiation in some markets is measured by the ability to provide reliable service in times of volatility rather than optimized costs in times of stability. Agentic supply chains enable continuous rebalancing of inventory; adjustments to routing; and renegotiation of sourcing agreements within predetermined policy boundaries; thereby providing companies with the opportunity to provide superior service fulfillment in comparison to competitors who experience stockouts; delays; or compliance bottlenecks (Christopher & Peck, 2004). Superior service fulfillment translates into customer loyalty; higher retention rates; and stronger brand preferences; which become strategically valuable assets beyond the traditional operational metrics.

The strategic value of governing autonomy is also reflected in how organizations allocate human managerial attention. Human managers have historically been required to focus on operational issues in order to ensure successful execution of day-to-day operations. However, with agentic systems managing routine decision-execution within predetermined governance constraints; human managers can refocus their attention on strategic design of supply networks; strategic design of supplier portfolios; and strategic design of risk policies (Ketchen & Hult, 2007). Redirection of managerial attention is economically significant due to the scarcity of managerial attention and its propensity to be wasted on variability in execution. Governance of autonomy provides organizations with the opportunity to reduce variability in decision making and decrease the amount of exceptions that require managerial intervention; thereby increasing an organization's bandwidth for innovation; development of suppliers; and partnership development. The supply chain can then serve as a strategic experimentation platform rather than simply a reactive function.

Competitive advantage through governing autonomy is dependent upon the credibility of governance mechanisms that allow stakeholders to have faith in the decisions being made. In regulated industries and global trade networks; stakeholders (e.g., partners; government agencies) evaluate how decisions are made and if compliance is integrated systematically. Companies that demonstrate audit-readiness; decision-traceability; and continuous compliance-monitoring; are able to receive faster approval; less friction; and greater confidence from partners (Kache & Seuring, 2017). Credibility of governance mechanisms can compound benefits over time since trusted autonomy enables greater integration with partners and automation of cross-enterprise coordination; thereby creating network effects that competitors may find difficult to replicate.

Another area of strategic implication for agentic supply chains relates to resilience and adaptability, where agentic supply chains transition resilience from manual contingency planning to continuous adaptive control (Tang, 2006). Traditional resilience strategies rely on redundancy buffers and human-led response teams. However, agentic systems provide a new level of resilience by developing the capability to anticipate disruptions

and dynamically reconfigure flows based on real-time signals. Strategic value arises from minimizing the time between a disruption emerges and the organization responds to the disruption; thereby preventing cascading failures that would otherwise negatively affect an organization's revenue and customer trust (Christopher & Peck, 2004). Therefore, resilience becomes an operational capability that directly impacts an organization's competitiveness.

In addition to the strategic value of adaptive resilience, adaptive resilience also provides an additional strategic benefit by allowing organizations to optimize efficiency while simultaneously optimizing robustness. Most conventional resilience strategies result in increased cost through inventory buffers or redundant capacity. Agentic systems can achieve similar levels of resilience through intelligent reallocation and targeted redundancy activation; thereby decreasing the need for blanket buffering. This capability improves an organization's capital efficiency and working capital management; thereby enabling an organization to sustain competitive pricing and profitability in volatile environments.

Finally, there are several implications for strategic risk governance and enterprise valuation. Investors and boards of directors increasingly evaluate an organization's operational resilience as an indicator of an organization's overall financial stability. Companies that demonstrate lower disruption-impact and faster recovery may receive lower-risk premiums and/or premium valuations (Teece, 2007). Agentic supply chains that measure and report resilience metrics; recovery performance; and governance-based evidence provide credible narratives that support such claims. Therefore, resilience is not only an operational concept, but also a financial and strategic concept that influences investors' perceptions of an organization's long-term viability.

Strategic flexibility under uncertainty is another strategic implication of agentic supply chains that reflects an organization's ability to adapt their supply chain configurations and operating policies as market conditions and regulatory requirements evolve. Sources of uncertainty in global supply chains include demand volatility; geopolitical shifts; trade restrictions; and data sovereignty constraints. Agentic systems provide strategic flexibility by enabling rapid policy updates; scenario evaluations; and controlled deployments of new decision-rules via governance-constrained mechanisms (Kache & Seuring, 2017). Strategic flexibility transitions from slow redesign cycles to continuous adaptation without sacrificing compliance or control.

Strategic flexibility also enables organizations to pursue new market opportunities with lower operational risk. Entry into new geographic markets often requires rapid adjustments to sourcing; logistics; and compliance practices. Agentic supply chains with governance-by-design can encode regional constraints; and learn local patterns rapidly while maintaining global oversight (Ketchen & Hult, 2007). Rapid market entry; and reduced ramp-up time; thereby enable organizations to pursue new market opportunities with greater strategic flexibility.

A conceptual model for describing flexibility can be represented through an options-value framework that describes the value derived from maintaining multiple feasible supply chain configurations that can be activated under different states of the world (Teece, 2007). If V represents expected value; and the action corresponds to configuration choices conditioned on states s , flexibility provides an organization with the opportunity to maximize value by aligning configuration with realized conditions. The strategic option logic that underlies the options-value framework helps explain why governing autonomy can increase long-term performance even though the short-term efficiencies provided by governing autonomy may be modest.

Competitive advantage, resilience, and flexibility are mutually reinforcing when autonomy is governed. Governing autonomy enables organizations to accelerate execution; reduce disruption-impacts; and provide strategic flexibility to reconfigure supply chain configurations under uncertainty. Collectively, these capabilities transform the supply chain from a tactical cost center into a strategic asset that influences customer experience; cost stability; compliance posture; and growth capacity (Tang, 2006). Therefore, firms that view agentic supply chains as governance-first systems, rather than solely as optimization engines, will develop more sustainable advantages.

Therefore, strategic and business implications of agentic supply chains exist beyond the realm of operational improvements into organizational transformations. Agentic supply chains influence decision authority; managerial roles; partner integration; and risk governance. Firms that invest in designing and implementing

governance architectures that incorporate mechanisms for auditability; resilience; and compliance-position themselves to compete in a future marked by disputed data environments; fragmented regulatory requirements; and increasing volatility (Waller & Fawcett, 2013). In this future; agentic supply chains can serve as a strategic differentiator-not because of the degree to which they automate tasks-but because of the degree to which they develop a governable; adaptive; and trustworthy decision infrastructure.

Ethical Implications of Delegated Decision Authority

Ethics of delegating decision-making authority to artificial agents (i.e., autonomous systems) is very important when supply chain control transitions from human-centered to autonomous agent-based systems (Floridi et al., 2018). Artificial agents that are capable of autonomous decision making change how moral responsibility, accountability, and social impacts are dispersed within organizations and society. Artificial agents will autonomously take action on the world by allocating resources, choosing suppliers, prioritizing customers, and determining what type of labor force is needed; these actions are ethically accountable and should not be treated as secondary concerns. Therefore, ethical legitimacy becomes a necessary condition for the long-term sustainability of autonomous supply chains.

A key challenge of designing autonomous supply chains is the need to redefine what constitutes "moral responsibility." As autonomous systems make decisions without a singular human actor, the decisions made by autonomous systems generate real-world consequences (Mittelstadt et al., 2016). For example, if an autonomous agent allocates its inventory to less vulnerable areas, delays a shipment of humanitarian aid, or chooses a supplier that has a questionable labor practice, the autonomous agent may cause some form of harm. However, it is impossible to assign moral responsibility to the artificial agent itself since moral agency is still a uniquely human construct. Instead, the organization that designed, configured, deployed and oversees the artificial agent is morally responsible for the autonomous agent's actions and behavior. Thus, ethical governance requires the explicit assignment of responsibility for the behavior of the autonomous system, as opposed to relying on a general sense of accountability.

One of the most critical challenges associated with the diffusion of responsibility is the possibility of moral disengagement (Floridi et al., 2018). If the organization believes that the autonomous agent, rather than the organization itself, has decision authority, then the organization may disengage itself from any moral responsibility for the adverse consequences resulting from the autonomous agent's actions. In supply chains, the risk of moral disengagement is compounded because many decisions involve difficult trade-offs between factors such as cost, service accessibility, and labor conditions. Therefore, ethical frameworks must ensure that the delegation of decision authority to autonomous systems does not diminish the organization's moral accountability, but instead, provides a new foundation for accountability through the establishment of governance roles, escalation mechanisms, and oversight mandates.

Another major ethical concern of autonomous supply chains is the potential for bias and unfairness. Autonomous decision systems are trained using data from history, which often reflect existing structural inequalities, regional disparities, and legacy relationships between suppliers and their customers (Jobin et al., 2019). If autonomous decision systems are allowed to operate unimpeded, they may exacerbate or continue to perpetuate these biases by consistently preferring specific suppliers, regions, or customer groups. The result could be the exclusion of small suppliers, marginalization of developing regions, or diminished service quality to lower margin customers. Thus, fairness cannot be assumed to emerge naturally from the optimization process, and fairness must be explicitly incorporated and monitored.

In addition, fairness risks are of particular concern in global supply chains, where power disparities between different economies are significant (Greene et al., 2019). Autonomous sourcing decisions may favor large suppliers with greater data availability over smaller suppliers who lack digital infrastructure. Although this concentration may improve efficiency, it may also undermine the long-term viability of small and medium-sized enterprises (SMEs), as well as the economic development of developing countries. Ethical assessment of autonomous decision systems therefore cannot focus exclusively on aggregate efficiency, but must also include consideration of the distributional effects of autonomous decisions.

In addition, bias in autonomous decision systems may arise from surrogate measures that are included in data, such as geographic region, delivery reliability, or historical performance metrics, which may correlate with socioeconomic characteristics (Mittelstadt et al., 2016). Autonomous systems may unintentionally discriminate while appearing to be neutral. Therefore, ethical oversight of autonomous systems requires ongoing audits of decision-making patterns to identify disparate treatment across regions, types of suppliers, or segments of the workforce. Auditing goes beyond compliance and includes a normative evaluation of whether outcomes align with the values of the organization and societal expectations.

Finally, workforce impacts are another ethical dimension of delegated decision authority. Autonomous supply chains modify job roles, decision authority, and required skills in logistics, procurement, and planning functions. Autonomous decision systems may eliminate certain jobs or downskill other jobs. However, ethical implementation of autonomous supply chains requires that organizations develop intentional work-force transition strategies, which prioritize reskilling, evolve job roles, and ensure meaningful human involvement in oversight and governance. Failure to address workforce impacts may lead to negative social reaction, loss of public trust, and reduced long-term stability for organizations.

The ethical implications of delegated decision authority extend beyond the internal workforce effects to broader societal outcomes. Supply chains affect employment patterns, environmental impact, access to essential goods and services, and regional economic stability (Owen et al., 2012). Autonomous decision systems that optimize narrowly for cost and speed may inadvertently harm communities through environmental degradation, labor exploitation, or denial of services. Therefore, autonomous supply chains must internalize societal impact considerations through governance constraints, performance metrics, and escalation mechanisms. This enables autonomy to be aligned with corporate social responsibility, rather than merely viewing ethics as an additional layer of regulation.

Responsible Innovation Principles Provide a Normative Framework for Designing Autonomous Supply Chains. Responsible innovation principles offer a normative framework for developing autonomous supply chains (Stilgoe et al., 2013). The principles emphasize anticipation, reflexivity, inclusion and responsiveness. Anticipation requires assessing potential ethical consequences of autonomy prior to deployment, including unforeseen second-order consequences. Reflexivity involves continuously assessing system behavior and ethical assumptions based on changing conditions. Inclusion requires collaboration with stakeholders who may be impacted by autonomous decision systems, including suppliers, workers, and communities. Responsiveness requires mechanisms to implement corrective actions when harm is identified.

By embedding responsible innovation principles in autonomous supply chains, ethics is transformed from a post-hoc review mechanism to a system-level design mechanism (Owen et al., 2012). Governance constraints, fairness objectives, workforce impacts, and societal risk thresholds become part of the decision logic, rather than simply being subject to external oversight. Integration of responsible innovation principles into autonomous supply chains reinforces ethical legitimacy, since the autonomous behavior is grounded in explicit moral commitments rather than implied technical priorities. Thus, ethical design is inherently connected to governance architecture.

Additionally, the ability to establish ethical legitimacy of autonomous supply chains positively influences strategic outcomes and adoption. Organizations that demonstrate responsible development and deployment of autonomous supply chains will build trust with regulators, partners, employees, and customers (Jobin et al., 2019). Trust leads to fewer obstacles to adoption, faster integration, and longer-term scalability. On the contrary, unethical behavior related to autonomous supply chains can create reputational damage, regulatory response, and social opposition, all of which can undermine technological advantages. Therefore, ethics has tangible strategic and economic consequences.

The ethical analysis of delegated decision authority can be framed conceptually by examining the expected ethical impact of a decision as a function of the decision probability and harm magnitude (Wachter et al., 2017). Using E to represent expected ethical impact, H to represent harm severity, and P to represent decision probability, governance attempts to minimize the expected harm through constraints, oversight, and escalation.

The conceptual framework illustrates that ethics risk management parallels operational risk management, except that ethics addresses normative outcomes rather than strictly financial outcomes.

Ultimately, the ethical implications of delegated decision authority will determine whether autonomous supply chains are socially acceptable and institutionally sustainable (Wachter et al., 2017). Therefore, organizations will ensure that autonomy contributes to the well-being of humans and society, rather than diminishing it. Ethical governance, therefore, is a prerequisite for the long-term trustworthiness, interdisciplinarity and socially responsible progress of autonomous supply chains.

Comparative Analysis with Existing Supply Chain Intelligence Models

A comparative analysis of agentic supply chains versus other supply chain intelligence models will help to clarify the degree to which agentic supply chains represent a conceptual novelty; the degree to which they offer an operational differentiation; and the degree to which they require changes in governance practices. Supply chain intelligence has developed progressively since the beginning of the 21st century via four successive paradigms: descriptive analytics, predictive forecasting, rule-based automation, and centralized visibility platforms (Waller & Fawcett, 2013). All of these paradigms have improved certain aspects of decision support. However, each paradigm retained a common structural assumption that humans remained the ultimate decision authority. In contrast, agentic supply chains depart from this assumption by providing autonomous decision-making capacity to act within delegated authority boundaries. This section will differentiate the proposed framework by describing how agentic decisioning fundamentally differs from predictive analytics and control tower architectures and rule-based systems regarding the extent of authority, flexibility, and governance (Choi et al., 2018). Predictive Analytics Paradigm

Predictive analytics is the predominant paradigm for intelligence in today's supply chains. Predominant predictive models predict demand, lead time, disruption probability, or future cost based upon past data and statistical learning (Gunasekaran et al., 2017). These models provide improved foresight but exist as advisory tools only. Humans analyze predictions make decisions and bear responsibility for the consequences of those decisions. This separation of functions provides humans with control over the decision-making process but creates latency, cognitive overload, and inconsistency. Although predictive analytics provides improved information quality, it fails to solve the problems associated with the frequent bottleneck of execution found in complex, global supply chains operating at high frequencies.

Agentic decision-making fundamentally differs from predictive analytics in its ability to close the gap between prediction and action. Learned policies within agentic systems select and execute actions such as redirecting inventory, re-routing shipments or negotiating sourcing agreements within predefined governance bounds. The key difference is not the improved predictive accuracy of agentic systems, but rather their ability to provide continuous, closed-loop control. Agentic systems respond to changing situations without waiting for humans to interpret the data. Therefore, predictive analytics transforms intelligence from foresight into behavior and from episodic human tasks into a system property.

Governance Perspective: Predictive analytics places accountability for outcomes clearly on human decision-makers, regardless of whether or not the outcome was influenced by the predictive analytics results. Conversely, agentic decision-making distributes accountability among system developers, governance stewards, and oversight structures. As a result, this redistribution requires the development of auditability, traceability, and constraint-based control mechanisms that are not required by predictive analytics models. The proposed framework addresses the identified gap by incorporating governance into the decision-making process at the time of decision-execution, rather than relying on post-hoc justification. This incorporation of governance represents a qualitative distinction from predictive analytics, rather than an incremental improvement.

Control Towers Paradigm: Control towers are another widely adopted paradigm for supply chain intelligence. Control towers are centralized repositories of data collected from across the entire supply network that provide real-time visibility, alerts, and dashboards (Barreto et al., 2017). The primary purpose of control towers is situational awareness, rather than decision execution. Advanced control towers may provide recommendations to take specific actions; however, they generally depend upon human operators to approve and execute those

recommendations. Therefore, control towers enhance coordination and transparency but fail to eliminate the human-mediated bottleneck and lack consistency in response across decision domains.

Digital Twin Orchestration – Beyond Visibility: Digital twin orchestration within agentic supply chains expands the scope of control towers' focus on visibility. Digital twins in the proposed framework are active execution substrates that facilitate autonomous decisions through real-time state synchronization and simulation (Ivanov & Dolgui, 2020). Unlike control towers that display the state of the supply chain, digital twins enact decisions by determining if an action is physically possible and compliant with regulatory requirements before taking the action. This operational role transforms digital twins from representative artifacts into control infrastructure. The difference is fundamental, as it permits safe autonomy instead of merely enhanced monitoring.

Scale and Responsiveness Limitations: Both control towers and predictive analytics are limited by scale and responsiveness as the complexity of the supply chain increases. Human operators become overwhelmed by the volume of alerts and exceptions, resulting in delayed or inconsistent responses to changing conditions. Agentic systems that operate through digital twin orchestration enable the distribution of decision-execution while maintaining centralized governance. This architecture enables scalability without compromising control. The comparative advantage of agentic systems lies not in richer dashboards but in eliminating the human bottleneck in routine decisions, while retaining oversight authority.

Rule-Based Systems Paradigm: Rule-based systems are an earlier automation paradigm in supply chain management. These systems codify deterministic logic such as reorder points, routing rules, and supplier selection criteria (Tang, 2006). Rule-based systems improve consistency and speed of execution but lack adaptability. Rule-based systems may generate suboptimal or even hazardous outcomes when conditions vary from the predetermined scenarios. Maintaining rule-sets grows increasingly difficult as supply chains globalize and regulations multiply.

Differences Between Learning Based Agentic Systems and Rule-Based Systems: Learning-based agentic systems differ from rule-based systems by adapting policies based upon experience rather than static logic. Learning-based systems generalize across scenarios and adjust their behavior as conditions evolve. However, learning alone does not guarantee safety or alignment. Without governance, learning-based systems may drift, exploit loopholes, or prioritize short-term objectives. The proposed framework distinguishes itself by constraining learning within rule-bounded policy spaces that enforce governance architecture. This combination maintains adaptability while ensuring that learning-based systems do not exhibit uncontrolled behavior.

Comparative Distinction: The comparative distinction between rule-based and learning-based systems can be described using policy selection dynamics. Deterministic action selection occurs in rule-based systems given a state. Agentic systems determine action selection by maximizing expected utility subject to constraints. If U denotes expected utility and P denotes policy selection over states s and actions a then agentic decisioning selects actions to maximize:

$$E(U) = \sum_s P(s | a) R(s, a)$$

within governance constraints that restrict admissible actions. This formulation illustrates that learning-based autonomy optimizes behavior within predefined boundaries rather than executing deterministically defined logic.

Treating Governance as External Layers: All existing intelligence models treat governance as an external layer that is applied after decisions are made. Predictive analytics relies on human judgment, control towers rely on human operator intervention, and rule-based systems rely on static constraints. The proposed agentic framework embeds governance into the decision-making process itself. Governance constraints, audit logging, escalation thresholds, and rollback mechanisms are evaluated at the time of decision-execution rather than retrospectively. This integration is the fundamental differentiator that enables scalable trust.

Business Advantages of Agentic Supply Chains: The comparative advantages of agentic supply chains include sustained performance in volatile conditions. Predictive analytics and control towers are effective in stable conditions but degrade rapidly under conditions of rapid change due to human bottlenecks. Rule-based systems fail under novel conditions. Governed-agentic systems maintain consistent and responsive decision-making and compliance under conditions of rapid change. This characteristic supports sustainable strategic resilience rather than episodic improvements in efficiency (Xu et al., 2018). The comparative analysis demonstrates that the proposed framework is not a rebranding of existing frameworks but a structural reorganization of supply chain intelligence. By transforming from advisory analytics and visibility platforms to governed-autonomous execution mediated by digital twins, the proposed framework establishes a new class of supply chain systems. This new class of systems addresses the deficiencies of prior paradigms while preserving accountability and control.

Future Research Directions

The future of agentic supply chain research should take a holistic view toward recognizing the increasing importance of autonomous decision systems as enduring organizational actors, as opposed to temporary technological tools. Because autonomous decision systems continue to make decisions regarding procurement logistics inventory management, and network coordination, the research community must begin to move away from individualized performance enhancements, and toward understanding how to govern the long-term systemic behavior of agentic supply chain systems, their sustainability, and their institutional integration. The continued relevance of this research agenda will depend upon its ability to answer persistent questions as the maturity of technologies continues, regulatory environments continue to evolve, and organizations continue to increase their dependence on autonomous systems.

One of the most significant and enduring research areas for agentic supply chains is multi-agent coordination. Supply chains are inherently multi-actor systems; the decisions made by one actor impact the operating environment of all other actors through shared resources, constraints, and feedback loops. In an agentic environment, coordination is no longer facilitated solely through centralized planning or human negotiation, but through the interaction of autonomous policies. Future research must identify the conditions under which local optimizations lead to global stability, and when they lead to systemic oscillations.

In addition to identifying the conditions under which coordination occurs, future research must also consider scale and heterogeneity. As the number of agents increases, so too does the complexity of coordination, leading to non-linear increases in the likelihood of congestion, amplification of inventory oscillations, and decision deadlocks. Agents will have differing learning models, data access, governance constraints, and risk tolerances. Each of these factors will influence the interaction dynamics and create potential for asymmetric power distributions throughout the system. Therefore, research that creates models of heterogeneous populations of agents, and identifies coordination mechanisms that are robust to diversity, will be required for real-world deployments.

Research into cross-enterprise agentic supply chains will further complicate the issues associated with coordination, while creating new challenges and opportunities for understanding issues of trust, legitimacy, and accountability. Autonomous agents representing different enterprises will have decisions that affect shared infrastructure, contractual obligations, and regulatory exposures. There will be no single entity capable of imposing unilateral governance. Future research must investigate decentralized governance mechanisms that facilitate cooperation without centralized authority. These mechanisms may incorporate elements of institutional theory, contract theory, and distributed systems research to provide common norms, enforcement protocols, and dispute resolution procedures. Additionally, cross-enterprise autonomy raises fundamental concerns regarding data sharing and competitive sensitivities. Enterprises may be reluctant or legally prohibited from sharing detailed operational data; however, they may still require cooperative decision-making. Therefore, research into privacy-preserving coordination techniques (e.g., constrained information exchange, abstracted state representations, and collaborative policy alignment) is essential. These techniques must balance the need for coordination efficiency with the requirement to protect proprietary information.

Another emerging area of research related to agentic supply chains is the use of quantum and hybrid optimization paradigms to improve the way agentic supply chains make complex decisions. A significant portion of the challenges associated with supply chain design, including network design, routing, and capacity allocation, exhibit combinatorial complexity that limits the effectiveness of traditional classical methods. The emerging computational paradigms may provide alternative approaches for exploring large solution spaces and generating high-quality candidate solutions. However, future research should focus on evaluating the practical advantages of these paradigms in solving complex decision problems, rather than simply speculating about potential performance benefits. The integration of advanced optimization paradigms into agentic systems raises several additional governance and interpretability challenges that must be addressed through scholarship. For example, novel solvers may generate output that contains unfamiliar types of uncertainty, or provide little transparency regarding the generation process. Research is needed to understand how to validate, constrain, and audit such outputs, within governance-first architectures. Potential hybrids of exploratory computation and conservative validation may represent practical compromises for meeting the needs of both exploration and validation. Finally, regulatory co-evolution represents a core and enduring research area due to the fact that autonomous supply chains challenge many of the current assumptions embedded in legal frameworks, including those related to human intent, episodic decision-making, and ex-post enforcement. Regulatory institutions currently operate under the assumption that autonomous systems operate intermittently, not continuously; and that they do not have the same adaptive capabilities as humans. Therefore, future research must explore how regulatory institutions will adapt to these new realities, and how agentic systems will be designed to anticipate and support the evolving regulatory expectations of regulatory institutions. Examples of this include ongoing compliance monitoring, embedded auditability, and mechanisms for dynamic reporting. Co-evolutionary research should also study the feedback loops between regulatory development and technological development. As regulatory bodies develop new requirements for algorithmic accountability, transparency, and human oversight, system architectures will adapt to meet these new expectations, potentially allowing for even more sophisticated regulatory approaches based on real-time monitoring and assessment, rather than periodic audits.

REFERENCES:

1. Abbas, A. E., van Velzen, T., Ofe, H., van de Kaa, G., Zuiderwijk, A., & de Reuver, M. (2024). Beyond control over data: Conceptualizing data sovereignty from a social contract perspective. *Electronic Markets*, 34, Article 20. <https://doi.org/10.1007/s12525-024-00695-2>
2. Abideen, A. Z., Sundram, V. P. K., Pyeman, J., Othman, A. K., & Sorooshian, S. (2021). Digital twin integrated reinforced learning in supply chain and logistics. *Logistics*, 5(4), 84. <https://doi.org/10.3390/logistics5040084>
3. Abouelrous, A., Faury, O., & Masson, R. (2023). Digital twin applications in urban logistics: An overview. *Transportmetrica A: Transport Science*. <https://doi.org/10.1080/21650020.2023.2216768>
4. Alles, M. G., Kogan, A., & Vasarhelyi, M. A. (2008). Putting continuous auditing theory into practice: Lessons from two pilot implementations. *Journal of Information Systems*, 22(2), 195–214. <https://doi.org/10.2308/jis.2008.22.2.195>
5. AlMulhim, A. F. (2021). Smart supply chain and firm performance: The role of digital technologies. *Business Process Management Journal*, 27(5), 1353–1372. <https://doi.org/10.1108/BPMJ-12-2020-0573>
6. Alongi, A., Giallorenzo, S., Lanese, I., & Mauro, J. (2022). An event sourced and observable architecture for microservice based systems. *Software: Practice and Experience*, 52(8), 1712–1739. <https://doi.org/10.1002/spe.3116>
7. Alshiekh, M., Bloem, R., Ehlers, R., Könighofer, B., Niekum, S., & Topcu, U. (2018). Safe reinforcement learning via shielding. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), 2669–2678. <https://doi.org/10.1609/aaai.v32i1.11797>
8. Altman, E. (1996). Constrained Markov decision processes with total cost criteria: Occupation measures and primal LP. *Mathematical Methods of Operations Research*, 43(1), 45–72. <https://doi.org/10.1007/BF01303434>
9. Altman, E. (1999). *Constrained Markov decision processes*. Chapman & Hall/CRC. <https://doi.org/10.1201/9781315140223>
10. Amato, C. (2024). An introduction to centralized training for decentralized execution in cooperative multi agent reinforcement learning. *arXiv*. <https://doi.org/10.48550/arXiv.2409.03052>

11. Atreyi Kankanhalli, Hua (Jonathan) Ye, Hock Hai Teo; Comparing Potential and Actual Innovators: An Empirical Study of Mobile Data Services Innovation1. *MIS Quarterly* 1 September 2015; 39 (3): 667–682. <https://doi.org/10.25300/MISQ/2015/39.3.07>
12. Barreto, L., Amaral, A., & Pereira, T. (2017). Industry 4.0 implications in logistics: An overview. *Procedia Manufacturing*, 13, 1245–1252. <https://doi.org/10.1016/j.promfg.2017.09.045>
13. Barykin, S. Y., Bochkarev, A. A., Kalina, O. V., & Yadykin, V. K. (2020). Concept for a supply chain digital twin. *International Journal of Mathematical, Engineering and Management Sciences*, 5(6), 1498–1515.
14. Batini, C., Cappiello, C., Francalanci, C., & Maurino, A. (2009). Methodologies for data quality assessment and improvement. *ACM Computing Surveys*, 41(3), Article 16. <https://doi.org/10.1145/1541880.1541883>
15. Benlian, A., Hess, T., & Buxmann, P. (2009). Drivers of SaaS adoption: An empirical study of different application types. *Business and Information Systems Engineering*, 1(5), 357–369. <https://doi.org/10.1007/s12599-009-0068-x>
16. Bernstein, D. S., Givan, R., Immerman, N., & Zilberstein, S. (2002). The complexity of decentralized control of Markov decision processes. *Mathematics of Operations Research*, 27(4), 819–840. <https://doi.org/10.1287/moor.27.4.819.297>
17. Bifet, A., & Gavalda, R. (2007). Learning from time changing data with adaptive windowing. In *Proceedings of the 2007 SIAM International Conference on Data Mining* (pp. 443–448). <https://doi.org/10.1137/1.9781611972771.42>
18. Böhmecke-Schwafert, M. (2024). The role of auditability in AI governance: Evidence and implications for regulation. *Telecommunications Policy*, 48(8), 102835. <https://doi.org/10.1016/j.telpol.2024.102835>
19. Bonawitz, K., Ivanov, V., Kreuter, B., Marcedone, A., McMahan, H. B., Patel, S., ... Seth, K. (2017). Practical secure aggregation for privacy preserving machine learning. *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 1175–1191. <https://doi.org/10.1145/3133956.3133982>
20. Bonomi, F., Milito, R., Zhu, J., & Addepalli, S. (2012). Fog computing and its role in the internet of things. *Proceedings of the First Edition of the MCC Workshop on Mobile Cloud Computing*, 13–16. <https://doi.org/10.1145/2342509.2342513>
21. Bose, R. P. J. C., van der Aalst, W. M. P., Zliobaite, I., & Pechenizkiy, M. (2014). Dealing with concept drift in process mining. *IEEE Transactions on Neural Networks and Learning Systems*, 25(1), 154–171. <https://doi.org/10.1109/TNNLS.2013.2278313>
22. Brailsford, S. C., Harper, P. R., Patel, B., & Pitt, M. (2009). An analysis of the academic literature on simulation and modelling in healthcare. *Journal of Simulation*, 3(3), 130–140. <https://doi.org/10.1057/jos.2009.10>
23. Brodersen, K. H., Gallusser, F., Koehler, J., Remy, N., & Scott, S. L. (2015). Inferring causal impact using Bayesian structural time series models. *The Annals of Applied Statistics*, 9(1), 247–274. <https://doi.org/10.1214/14-AOAS788>
24. Brunke, L., Greeff, M., Hall, A. W., Yuan, Z., Zhou, S., Panerati, J., & Schoellig, A. P. (2022). Safe learning in robotics: From learning based control to safe reinforcement learning. *Annual Review of Control, Robotics, and Autonomous Systems*, 5, 411–444. <https://doi.org/10.1146/annurev-control-042920-020211>
25. Burgos, D., & Ivanov, D. (2021). Food retail supply chain resilience and the COVID-19 pandemic: A digital twin based impact analysis and improvement directions. *Transportation Research Part E: Logistics and Transportation Review*, 152, 102412. <https://doi.org/10.1016/j.tre.2021.102412>
26. Busoniu, L., Babuska, R., & De Schutter, B. (2008). A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics Part C*, 38(2), 156–172. <https://doi.org/10.1109/TSMCC.2007.913919>
27. Busse, A., Gerlach, B., Lengeling, J. C., Poschmann, P., Werner, J., & Zarnitz, S. (2021). Towards Digital Twins of Multimodal Supply Chains. *Logistics*, 5(2), 25. <https://doi.org/10.3390/logistics5020025>
28. Butner K (2010), "The smarter supply chain of the future". *Strategy & Leadership*, Vol. 38 No. 1 pp. 22–31, doi: <https://doi.org/10.1108/10878571011009859>
29. Carlini, N., & Wagner, D. (2017). Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)* (pp. 39–57). <https://doi.org/10.1109/SP.2017.49>

30. Catalano, M., Chirurgo, A., Fusto, C., Gazzaneo, L., Longo, F., Mirabelli, G., Nicoletti, L., Solina, V., & Talarico, S. (2022). A digital twin driven and conceptual framework for enabling extended reality applications: A case study of a brake discs manufacturer. *Procedia Computer Science*, 200, 1885–1893. <https://doi.org/10.1016/j.procs.2022.01.389>
31. Chaharsooghi, S. K., Heydari, J., & Zegordi, S. H. (2008). A reinforcement learning model for supply chain ordering management: An application to the beer game. *Decision Support Systems*, 45(4), 949–959. <https://doi.org/10.1016/j.dss.2008.03.007>
32. Chalendar, C., Raballand, G., & Rakotoarisoa, A. (2019). The use of detailed statistical data in customs reform: Evidence on risk management and compliance. *Development Policy Review*, 37(S1), O197–O222. <https://doi.org/10.1111/dpr.12352>
33. Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A survey. *ACM Computing Surveys*, 41(3), Article 15. <https://doi.org/10.1145/1541880.1541882>
34. Cheney, J., Chiticariu, L., & Tan, W. C. (2009). Provenance in databases: Why, how, and where. *Foundations and Trends in Databases*, 1(4), 379–474. <https://doi.org/10.1561/1900000006>
35. Cheong, B. C. C. (2024). Transparency and accountability in AI systems: Safeguarding wellbeing in the age of algorithmic decision making. *Frontiers in Human Dynamics*, 6, 1421273. <https://doi.org/10.3389/fhumd.2024.1421273>
36. Choi, T. M., Wallace, S. W., & Wang, Y. (2018). Big data analytics in operations management. *Production and Operations Management*, 27(10), 1868–1883. <https://doi.org/10.1111/poms.12838>
37. Chow, Y., Nachum, O., Duenez Guzman, E., & Ghavamzadeh, M. (2018). A Lyapunov based approach to safe reinforcement learning. *Advances in Neural Information Processing Systems*, 31, 8103–8112. <https://doi.org/10.48550/arXiv.1805.07708>
38. Christopher, M., & Peck, H. (2004). Building the resilient supply chain. *The International Journal of Logistics Management*, 15(2), 1–13. <https://doi.org/10.1108/09574090410700275>
39. Dai, J., & Vasarhelyi, M. A. (2017). Toward blockchain-based accounting and assurance: Continuous audit implications. *Journal of Information Systems*, 31(3), 5–21. <https://doi.org/10.2308/isys-51804>
40. Dietterich, T. G. (2000). Ensemble methods in machine learning. In *Multiple Classifier Systems* (pp. 1–15). https://doi.org/10.1007/3-540-45014-9_1
41. Dominguez, R., & Cannella, S. (2020). Insights on multi agent systems applications for supply chain management. *Sustainability*, 12(5), 1935. <https://doi.org/10.3390/su12051935>
42. Doshi Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *arXiv*. <https://doi.org/10.48550/arXiv.1702.08608>
43. Dzindolet, M. T., Peterson, S. A., Pomranky, R. A., Pierce, L. G., & Beck, H. P. (2003). The role of trust in automation reliance. *International Journal of Human Computer Studies*, 58(6), 697–718. [https://doi.org/10.1016/S1071-5819\(03\)00038-7](https://doi.org/10.1016/S1071-5819(03)00038-7)
44. Endsley, M. R., & Kaber, D. B. (1999). Level of automation effects on performance, situation awareness and workload in a dynamic control task. *Ergonomics*, 42(3), 462–492. <https://doi.org/10.1080/001401399185595>
45. Feinberg, E. A., & Schwartz, A. (1995). Constrained Markov decision models with weighted discounted rewards. *Mathematics of Operations Research*, 20(2), 302–320. <https://doi.org/10.1287/moor.20.2.302>
46. Floridi, L., & Cowls, J. (2019). A unified framework of five principles for AI in society. *Harvard Data Science Review*, 1(1). <https://doi.org/10.1162/99608f92.8cd550d1>
47. Floridi, L., Cowls, J., Beltramini, M., Chatila, R., Chazerand, P., Dignum, V., Luetge, C., Madelin, R., Pagallo, U., Rossi, F., Schafer, B., Valcke, P., & Vayena, E. (2018). AI4People—An ethical framework for a good AI society. *Minds and Machines*, 28, 689–707. <https://doi.org/10.1007/s11023-018-9482-5>
48. Frazzon, E. M., Agostino, I. R. S., Broda, E., & Freitag, M. (2020). Manufacturing networks in the era of digital production and operations: A socio cyber physical perspective. *Annual Reviews in Control*, 49, 288–294. <https://doi.org/10.1016/j.arcontrol.2020.04.008>
49. Fuller, A., Fan, Z., Day, C., & Barlow, C. (2020). Digital twin: Enabling technologies, challenges and open research. *IEEE Access*, 8, 108952–108971. <https://doi.org/10.1109/ACCESS.2020.2998358>
50. Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., & Bouchachia, A. (2014). A survey on concept drift adaptation. *ACM Computing Surveys*, 46(4), Article 44. <https://doi.org/10.1145/2523813>

51. Garvey, M. D., Carnovale, S., & Yenyurt, S. (2015). An analytical framework for supply network risk propagation: A Bayesian network approach. *European Journal of Operational Research*, 243(2), 618–627. <https://doi.org/10.1016/j.ejor.2014.10.034>
52. Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Daumé III, H., & Crawford, K. (2021). Datasheets for datasets. *Communications of the ACM*, 64(12), 86–92. <https://doi.org/10.1145/3458723>
53. Ghofrani, F., He, Q., Goverde, R. M. P., & Liu, X. (2018). Recent applications of big data analytics in railway transportation systems. *Transportation Research Part C*, 90, 226–246. <https://doi.org/10.1016/j.trc.2018.03.010>
54. Goodfellow, I. J., Shlens, J., & Szegedy, C. (2015). Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*. <https://doi.org/10.48550/arXiv.1412.6572>
55. Greene, D., Hoffmann, A. L., & Stark, L. (2019). Better, nicer, clearer, fairer: A critical assessment of the movement for ethical AI and fairness. *Proceedings of HICSS*, 1–10. <https://doi.org/10.24251/HICSS.2019.258>
56. Grieves, M., & Vickers, J. (2017). Digital twin: Mitigating unpredictable, undesirable emergent behavior in complex systems. In F. J. Kahlen, S. Flumerfelt, & A. Alves (Eds.), *Transdisciplinary perspectives on complex systems* (pp. 85–113). Springer. https://doi.org/10.1007/978-3-319-38756-7_4
57. Gu, S., Kelly, M., & others. (2024). A review of safe reinforcement learning: Methods, theories, and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12), 11216–11235. <https://doi.org/10.1109/TPAMI.2024.3457538>
58. Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2018). A survey of methods for explaining black box models. *ACM Computing Surveys*, 51(5), Article 93. <https://doi.org/10.1145/3236009>
59. Gunasekaran, A., Papadopoulos, T., Dubey, R., Fosso Wamba, S., Childe, S. J., Hazen, B., & Akter, S. (2017). Big data and predictive analytics for supply chain and organizational performance. *Journal of Business Research*, 70, 308–317. <https://doi.org/10.1016/j.jbusres.2016.08.004>
60. Guo, D., Zhong, R. Y., & Huang, G. Q. (2025). The role of digital twins in lean supply chain management. *International Journal of Production Research*. <https://doi.org/10.1080/00207543.2024.2372655>
61. Hancock, P. A., Billings, D. R., Schaefer, K. E., Chen, J. Y. C., de Visser, E. J., & Parasuraman, R. (2011). A meta analysis of factors affecting trust in human robot interaction. *Human Factors*, 53(5), 517–527. <https://doi.org/10.1177/0018720811417254>
62. Haripriya, R., Khare, N., Pandey, M., & Biswas, S. (2025). Navigating the fusion of federated learning and big data: A systematic review for the AI landscape. *Cluster Computing*, 28, 1–28. <https://doi.org/10.1007/s10586-024-05070-6>
63. Haskell, W. B., & Jain, R. (2013). Stochastic dominance constrained Markov decision processes. *SIAM Journal on Control and Optimization*, 51(1), 273–303. <https://doi.org/10.1137/120874679>
64. Haviv, M. (1996). On constrained Markov decision processes. *Operations Research Letters*, 19(1), 25–32. [https://doi.org/10.1016/0167-6377\(96\)00003-X](https://doi.org/10.1016/0167-6377(96)00003-X)
65. He, X., & Zhang, Y. (2023). Adversarial examples cybersecurity of deep learning: A survey of methods, applications, and challenges. *Expert Systems with Applications*, 230, 122223. <https://doi.org/10.1016/j.eswa.2023.122223>
66. Hellmeier, M., Pampus, J., Qarawlus, H., & Howar, F. (2023). Implementing data sovereignty: Requirements and challenges from practice. In *Proceedings of the 18th International Conference on Availability, Reliability and Security (ARES 2023)*. ACM. <https://doi.org/10.1145/3600160.3604995>
67. Heluany, J. B., et al. (2023). Survey on digital twins: From concepts to applications. *Proceedings of the ACM on Management of Data*. <https://doi.org/10.1145/3600160.3605070>
68. Hoff, K. A., & Bashir, M. (2015). Trust in automation: Integrating empirical evidence on factors that influence trust. *Human Factors*, 57(3), 407–434. <https://doi.org/10.1177/0018720814547570>
69. Hosseini, S., Ivanov, D., & Dolgui, A. (2019). Review of quantitative methods for supply chain resilience analysis. *Transportation Research Part E*, 125, 285–307. <https://doi.org/10.1016/j.tre.2019.03.001>
70. Huang, L., Joseph, A. D., Nelson, B., Rubinstein, B. I. P., & Tygar, J. D. (2011). Adversarial machine learning. In *Proceedings of the 4th ACM Workshop on Security and Artificial Intelligence* (pp. 43–58). <https://doi.org/10.1145/2046684.2046692>

71. Hummel, P., Braun, M., Tretter, M., & Dabrock, P. (2021). Data sovereignty: A review. *Big Data & Society*, 8(1). <https://doi.org/10.1177/2053951720982012>
72. Hunt, R., & Jackson, M. (2010). An introduction to continuous controls monitoring. *Computer Fraud & Security*, 2010(6), 16–19. [https://doi.org/10.1016/S1361-3723\(10\)70069-5](https://doi.org/10.1016/S1361-3723(10)70069-5)
73. Iftekhhar, A., Cui, X., Hassan, M. M., & Afzal, W. (2021). Blockchain-based traceability system that ensures food safety and quality. *Foods*, 10(6), 1289. <https://doi.org/10.3390/foods10061289>
74. Irfan, M., Malik, K., & Muhammad, K. (2024). Federated fusion learning with attention mechanism for multi client medical image analysis. *Information Fusion*, 108, 102364. <https://doi.org/10.1016/j.inffus.2024.102364>
75. Ivanov, D. (2017). Simulation based ripple effect modelling in the supply chain. *International Journal of Production Research*, 55(7), 2083–2101. <https://doi.org/10.1080/00207543.2016.1275873>
76. Ivanov, D. (2020). Predicting the impacts of epidemic outbreaks on global supply chains: A simulation based analysis on the coronavirus outbreak (COVID 19 SARS CoV 2) case. *Transportation Research Part E: Logistics and Transportation Review*, 136, 101922. <https://doi.org/10.1016/j.tre.2020.101922>
77. Ivanov, D. (2020). Viable supply chain model: Integrating agility, resilience and sustainability perspectives lessons from and thinking beyond the COVID 19 pandemic. *Annals of Operations Research*, 319, 1411–1431. <https://doi.org/10.1007/s10479-020-03640-6>
78. Ivanov, D., & Dolgui, A. (2020). A digital supply chain twin for managing disruption risks and resilience in the era of Industry 4.0. *Production Planning and Control*, 31(10), 775–788. <https://doi.org/10.1080/09537287.2020.1768450>
79. Jamshidi, P., Pahl, C., Mendonça, N. C., Lewis, J., & Tilkov, S. (2018). Microservices: The journey so far and challenges ahead. *IEEE Software*, 35(3), 24–35. <https://doi.org/10.1109/MS.2018.2141039>
80. Jannelli, V., Di Vaio, A., Palladino, R., & Schiraldi, M. M. (2025). Agentic LLMs in the supply chain: Towards autonomous decisioning. *International Journal of Production Research*. <https://doi.org/10.1080/00207543.2025.2604311>
81. Jans, M., Alles, M. G., & Vasarhelyi, M. A. (2014). A field study on the use of process mining of event logs as an analytical procedure in auditing. *The Accounting Review*, 89(5), 1751–1773. <https://doi.org/10.2308/accr-50807>
82. Jobin, A., Ienca, M., & Vayena, E. (2019). The global landscape of AI ethics guidelines. *Nature Machine Intelligence*, 1, 389–399. <https://doi.org/10.1038/s42256-019-0088-2>
83. Kaber, D. B. (2018). Issues in human automation interaction modeling: Presumptive aspects of frameworks of types and levels of automation. *Journal of Cognitive Engineering and Decision Making*, 12(1), 7–24. <https://doi.org/10.1177/1555343417737203>
84. Kaber, D. B., & Endsley, M. R. (1997). Out of the loop performance problems and the use of intermediate levels of automation for improved control system functioning and safety. *Process Safety Progress*, 16(3), 126–131. <https://doi.org/10.1002/prs.680160304>
85. Kache, F., & Seuring, S. (2017). Challenges and opportunities of digital information at the intersection of big data analytics and supply chain management. *International Journal of Operations and Production Management*, 37(1), 10–36. <https://doi.org/10.1108/IJOPM-02-2015-0078>
86. Kacianka, S., & Pretschner, A. (2021). Designing accountable systems. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. ACM. <https://doi.org/10.1145/3442188.3445905>
87. Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., ... Zhao, S. (2021). Advances and open problems in federated learning. *Foundations and Trends in Machine Learning*, 14(1–2), 1–210. <https://doi.org/10.1561/22000000083>
88. Ketchen, D. J., & Hult, G. T. M. (2007). Bridging organization theory and supply chain management. *Journal of Operations Management*, 25(2), 573–580. <https://doi.org/10.1016/j.jom.2006.05.010>
89. Kim, B., Kim, J. G., & Lee, S. (2024). A multi agent reinforcement learning model for inventory transshipments under supply chain disruption. *IIEE Transactions*, 56(7), 715–728. <https://doi.org/10.1080/24725854.2023.2217248>
90. Klein, G., Woods, D. D., Bradshaw, J. M., Hoffman, R. R., & Feltovich, P. J. (2004). Ten challenges for making automation a “team player” in joint human agent activity. *IEEE Intelligent Systems*, 19(6), 91–95. <https://doi.org/10.1109/MIS.2004.74>

91. Koot, M., Mes, M. R. K., & Iacob, M. E. (2021). A systematic literature review of supply chain digital twins. *Computers & Industrial Engineering*, 157, 107076. <https://doi.org/10.1016/j.cie.2020.107076>
92. Kritzinger, W., Karner, M., Traar, G., Henjes, J., & Sihn, W. (2018). Digital twin in manufacturing: A categorical literature review and classification. *IFAC PapersOnLine*, 51(11), 1016–1022. <https://doi.org/10.1016/j.ifacol.2018.08.474>
93. Kshetri, N. (2017). Blockchain's roles in strengthening cybersecurity and protecting privacy. *Telecommunications Policy*, 41(10), 1027–1038. <https://doi.org/10.1016/j.telpol.2017.09.003>
94. Kshetri, N. (2018). 1 Blockchain's roles in meeting key supply chain management objectives. *International Journal of Information Management*, 39, 80–89. <https://doi.org/10.1016/j.ijinfomgt.2017.12.005>
95. Kuehn, W. (2018). Digital twins for decision making in complex production and logistic enterprises. *International Journal of Design and Nature and Ecodynamics*, 13(3), 260–271. <https://doi.org/10.2495/DNE-V13-N3-260-271>
96. Kushwaha, A., Ravish, M., & others. (2025). A survey of safe reinforcement learning and constrained Markov decision processes. *arXiv*. <https://doi.org/10.48550/arXiv.2505.17342>
97. Lazarus, C., Lopez, J., & Kochenderfer, M. J. (2020). Runtime safety assurance using reinforcement learning. In 2020 IEEE AIAA 39th Digital Avionics Systems Conference (DASC). <https://doi.org/10.1109/DASC50938.2020.9256446>
98. LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444. <https://doi.org/10.1038/nature14539>
99. Lee, I., & Lee, K. (2015). The internet of things (IoT): Applications, investments, and challenges. *Business Horizons*, 58(4), 431–440. <https://doi.org/10.1016/j.bushor.2015.03.008>
100. Lee, J. H., Kim, C. O., & Park, S. J. (2008). Multi agent systems applications in manufacturing systems and supply chain management: A review paper. *International Journal of Production Research*, 46(1), 233–265. <https://doi.org/10.1080/00207540701441921>
101. Lee, J., Bagheri, B., & Kao, H. A. (2015). A cyber physical systems architecture for Industry 4.0 based manufacturing systems. *Manufacturing Letters*, 3, 18–23. <https://doi.org/10.1016/j.mfglet.2014.12.001>
102. Li, X., Huang, K., Yang, W., Wang, S., & Zhang, Z. (2020). On the convergence of FedAvg on non IID data. *Proceedings of ICLR 2020 Workshop*. <https://doi.org/10.48550/arXiv.1907.02189>
103. Li, Y., & Goel, S. (2025). Bridging IT auditors and AI auditing: Understanding pathways to effective IT audits of AI driven processes. *Advances in Accounting*, 69, 100842. <https://doi.org/10.1016/j.adiac.2025.100842>
104. Li, Y., Zobel, C. W., Seref, O., & Chatfield, D. (2020). Network characteristics and supply chain resilience under conditions of risk propagation. *International Journal of Production Economics*, 223, 107529. <https://doi.org/10.1016/j.ijpe.2019.107529>
105. Lipton, Z. C. (2018). The mythos of model interpretability. *Communications of the ACM*, 61(10), 36–43. <https://doi.org/10.1145/3233231>
106. Madhavan, P., & Wiegmann, D. A. (2007). Similarities and differences between human human and human automation trust. *Theoretical Issues in Ergonomics Science*, 8(4), 277–301. <https://doi.org/10.1080/14639220500337708>
107. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., & Vladu, A. (2018). Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*. <https://doi.org/10.48550/arXiv.1706.06083>
108. Marmolejo Saucedo, J. A. (2020). Design and development of digital twins: A case study in supply chains. *Mobile Networks and Applications*, 25(6), 2141–2160. <https://doi.org/10.1007/s11036-020-01557-9>
109. McMahan, H. B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication efficient learning of deep networks from decentralized data. *Proceedings of AISTATS 2017 (PMLR 54)*, 1273–1282. <https://doi.org/10.48550/arXiv.1602.05629>
110. Merritt, S. M., Heimbaugh, H., LaChapell, J., & Lee, D. (2013). I trust it, but I don't know why: Effects of implicit attitudes toward automation on trust in automation. *Human Factors*, 55(3), 520–534. <https://doi.org/10.1177/0018720812465081>

111. Mirsky, Y., Doitshman, T., Elovici, Y., & Shabtai, A. (2018). Kitsune: An ensemble of autoencoders for online network intrusion detection. In *Network and Distributed System Security Symposium (NDSS)*. <https://doi.org/10.14722/ndss.2018.23204>
112. Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., Spitzer, E., Raji, I. D., & Gebru, T. (2019). Model cards for model reporting. *Proceedings of the Conference on Fairness, Accountability, and Transparency (FAT '19)**, 220–229. <https://doi.org/10.1145/3287560.3287596>
113. Mittelstadt, B. D., Allo, P., Taddeo, M., Wachter, S., & Floridi, L. (2016). The ethics of algorithms: Mapping the debate. *Big Data and Society*, 3(2), 1–21. <https://doi.org/10.1177/2053951716679679>
114. Monostori, L., Kádár, B., Bauernhansl, T., Kondoh, S., Kumara, S., Reinhart, G., Sauer, O., Schuh, G., Sihn, W., & Ueda, K. (2016). Cyber physical systems in manufacturing. *CIRP Annals*, 65(2), 621–641. <https://doi.org/10.1016/j.cirp.2016.06.005>
115. Moreau, L., et al. (2015). The rationale of PROV. *Web Semantics: Science, Services and Agents on the World Wide Web*, 35, 235–257. <https://doi.org/10.1016/j.websem.2015.04.001>
116. Mousa, M., van de Berg, D., Kotecha, N., del Rio Chanona, E. A., & Mowbray, M. (2024). An analysis of multi agent reinforcement learning for decentralized inventory control systems. *Computers & Chemical Engineering*, 186, 108783. <https://doi.org/10.1016/j.compchemeng.2024.108783>
117. Nakao, K., Conroy, K., & Wen, Z. (2021). Data robust partially observable Markov decision processes. *SIAM Journal on Optimization*, 31(4), 2730–2757. <https://doi.org/10.1137/19M1268410>
118. Negri, E., Fumagalli, L., & Macchi, M. (2017). A review of the roles of digital twin in CPS based production systems. *Procedia Manufacturing*, 11, 939–948. <https://doi.org/10.1016/j.promfg.2017.07.198>
119. Nikolay Archak, Anindya Ghose, Panagiotis G. Ipeirotis, (2011) Deriving the Pricing Power of Product Features by Mining Consumer Reviews. *Management Science* 57(8):1485-1509.
120. Norman, D. A. (1990). The problem of automation: Inappropriate feedback and interaction, not “over automation”. *Philosophical Transactions of the Royal Society B*, 327(1241), 585–593. <https://doi.org/10.1098/rstb.1990.0101>
121. Novelli, C., Taddeo, M., & Floridi, L. (2024). Accountability in artificial intelligence: What it is and how it works. *AI & Society*, 39, 1871–1882. <https://doi.org/10.1007/s00146-023-01635-y>
122. Oroojlooyjadid, A., Nazari, M., Snyder, L. V., & Takáč, M. (2022). A deep Q network for the beer game: Deep reinforcement learning for inventory optimization. *Manufacturing & Service Operations Management*, 24(1), 285–304. <https://doi.org/10.1287/msom.2020.0939>
123. Owen, R., Macnaghten, P., & Stilgoe, J. (2012). Responsible research and innovation: From science in society to science for society, with society. *Science and Public Policy*, 39(6), 751–760. <https://doi.org/10.1093/scipol/scs093>
124. Panetto, H., Iung, B., Ivanov, D., Weichhart, G., & Wang, X. (2019). Challenges for the cyber physical manufacturing enterprises of the future. *Annual Reviews in Control*, 47, 200–213. <https://doi.org/10.1016/j.arcontrol.2019.02.002>
125. Papagiannidis, E., Mikalef, P., & Conboy, K. (2024). Responsible artificial intelligence governance: A review and research framework. *The Journal of Strategic Information Systems*, 33(4), 101885. <https://doi.org/10.1016/j.jsis.2024.101885>
126. Parasuraman, R., Sheridan, T. B., & Wickens, C. D. (2000). A model for types and levels of human interaction with automation. *IEEE Transactions on Systems, Man, and Cybernetics Part A: Systems and Humans*, 30(3), 286–297. <https://doi.org/10.1109/3468.844354>
127. Phiri, C. C. (2025). Creating characteristically auditable agentic AI systems. In *Proceedings of the Intelligent Robotics FAIR 2025 (IntRob '25)*. ACM. <https://doi.org/10.1145/3759355.3759356>
128. Piancastelli, C., & Tucci, M. (2020). The role of digital twins in the fulfilment logistics chain. *IFAC PapersOnLine*, 53(2), 10574–10578. <https://doi.org/10.1016/j.ifacol.2020.12.2807>
129. Raji, I. D., Smart, A., White, R. N., Mitchell, M., Gebru, T., Hutchinson, B., Smith Loud, J., Theron, D., & Barnes, P. (2020). Closing the AI accountability gap: Defining an end to end framework for internal algorithmic auditing. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (pp. 33–44). <https://doi.org/10.1145/3351095.3372873>
130. Reichert, M., & Weber, B. (2012). *Enabling flexibility in process aware information systems*. Springer.

131. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?" Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135–1144. <https://doi.org/10.1145/2939672.2939778>
132. Rodríguez Barroso, N., Stoyanov, S., & Gómez, J. (2023). Survey on federated learning threats: Concepts, taxonomy and defenses. *Information Fusion*, 92, 105–126. <https://doi.org/10.1016/j.inffus.2022.09.011>
133. Rolf, B., Jackson, I., Müller, M., Lang, S., Reggelin, T., & Ivanov, D. (2023). A review on reinforcement learning algorithms and applications in supply chain management. *International Journal of Production Research*, 61(20), 7151–7179. <https://doi.org/10.1080/00207543.2022.2140221>
134. Ross, K. W., & Varadarajan, R. (1989). Markov decision processes with sample path constraints: The communicating case. *Operations Research*, 37(5), 780–790. <https://doi.org/10.1287/opre.37.5.780>
135. Rozinat, A., & van der Aalst, W. M. P. (2008). Conformance checking of processes based on monitoring real behavior. *Information Systems*, 33(1), 64–95. <https://doi.org/10.1016/j.is.2007.07.001>
136. Saberi, S., Kouhizadeh, M., Sarkis, J., & Shen, L. (2019). Blockchain technology and its relationships to sustainable supply chain management. *International Journal of Production Research*, 57(7), 2117–2135. <https://doi.org/10.1080/00207543.2018.1533261>
137. Samuli Laato, Teemu Birkstedt, Matti Mäntymäki, Matti Minkkinen, and Tommi Mikkonen. 2022. AI governance in the system development life cycle: insights on responsible machine learning engineering. In *Proceedings of the 1st International Conference on AI Engineering: Software Engineering for AI (CAIN '22)*. Association for Computing Machinery, New York, NY, USA, 113–123. <https://doi.org/10.1145/3522664.3528598>
138. Sani, S., Zarifnia, A., Salontis, K., & Milisavljevic Syed, J. (2024). Supply Chain 4.0 and the digital twin approach: A framework for improving supply chain visibility. *Procedia CIRP*, 128, 321–326. <https://doi.org/10.1016/j.procir.2024.03.014>
139. Scerri, P., Pynadath, D. V., & Tambe, M. (2002). Towards adjustable autonomy for the real world. *Journal of Artificial Intelligence Research*, 17, 171–228. <https://doi.org/10.1613/jair.1037>
140. Schiff, D. S., Biddle, J., Borenstein, J., & Laas, K. (2024). The emergence of artificial intelligence ethics auditing. *Big Data & Society*, 11(2). <https://doi.org/10.1177/20539517241299732>
141. Schulman, J., Levine, S., Moritz, P., Jordan, M. I., & Abbeel, P. (2015). Trust region policy optimization. *arXiv*. <https://doi.org/10.48550/arXiv.1502.05477>
142. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv*. <https://doi.org/10.48550/arXiv.1707.06347>
143. Shapiro, A., Dentcheva, D., & Ruszczyński, A. (2014). *Lectures on stochastic programming: Modeling and theory* (2nd ed.). SIAM. <https://doi.org/10.1137/1.9781611973433>
144. Shneiderman, B. (2020). Human centered artificial intelligence: Reliable, safe and trustworthy. *International Journal of Human Computer Interaction*, 36(6), 495–504. <https://doi.org/10.1080/10447318.2020.1741118>
145. Simmhan, Y. L., Plale, B., & Gannon, D. (2005). A survey of data provenance in e science. *SIGMOD Record*, 34(3), 31–36. <https://doi.org/10.1145/1084805.1084812>
146. Stilgoe, J., Owen, R., & Macnaghten, P. (2013). Developing a framework for responsible innovation. *Research Policy*, 42(9), 1568–1580. <https://doi.org/10.1016/j.respol.2013.05.008>
147. Suriadi, S., Wynn, M. T., Xu, J., van der Aalst, W. M. P., & ter Hofstede, A. H. M. (2017). Event log imperfection patterns for process mining: Towards a systematic approach. *Information Systems*, 64, 132–150. <https://doi.org/10.1016/j.is.2016.07.011>
148. Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1), 9–44. <https://doi.org/10.1023/A:1022633531479>
149. Tang, C. S. (2006). Perspectives in supply chain risk management. *International Journal of Production Economics*, 103(2), 451–488. <https://doi.org/10.1016/j.ijpe.2005.12.006>
150. Tang, C. S. (2006). Robust strategies for mitigating supply chain disruptions. *International Journal of Logistics Research and Applications*, 9(1), 33–45. <https://doi.org/10.1080/13675560500405584>
151. Tao, F., Cheng, J., Qi, Q., Zhang, M., Zhang, H., & Sui, F. (2018). Digital twin driven product design, manufacturing and service with big data. *The International Journal of Advanced Manufacturing Technology*, 94(9–12), 3563–3576. <https://doi.org/10.1007/s00170-017-0233-1>

152. Tao, F., Qi, Q., Liu, A., & Kusiak, A. (2019). Digital twins and cyber physical systems toward smart manufacturing and Industry 4.0. *Engineering*, 5(4), 653–661. <https://doi.org/10.1016/j.eng.2019.01.014>
153. Tao, F., Zhang, H., Liu, A., & Nee, A. Y. C. (2019). Digital twin in industry: State of the art. *IEEE Transactions on Industrial Informatics*, 15(4), 2405–2415. <https://doi.org/10.1109/TII.2018.2873186>
154. Taylor, E. (2020). Data localization: The internet in the balance. *Telecommunications Policy*, 44(8), 102003. <https://doi.org/10.1016/j.telpol.2020.102003>
155. Teece, D. J. (2007). Explicating dynamic capabilities: The nature and microfoundations of enterprise performance. *Strategic Management Journal*, 28(13), 1319–1350. <https://doi.org/10.1002/smj.640>
156. Terrada, L., El Khaïli, M., & Ouajji, H. (2020). Multi agents system implementation for supply chain management making decision. *Procedia Computer Science*, 177, 624–630. <https://doi.org/10.1016/j.procs.2020.10.089>
157. Tong, X., Lai, K. H., Lo, C. K. Y., & Cheng, T. C. E. (2022). Supply chain security certification and operational performance: The role of upstream complexity. *International Journal of Production Economics*, 247, 108433. <https://doi.org/10.1016/j.ijpe.2022.108433>
158. Truong, C., Oudre, L., & Vayatis, N. (2020). Selective review of offline change point detection methods. *Signal Processing*, 167, 107299. <https://doi.org/10.1016/j.sigpro.2019.107299>
159. Uhlemann, T. H. J., Schock, C., Lehmann, C., Freiburger, S., & Steinhilper, R. (2017). The digital twin: Realizing the cyber physical production system for Industry 4.0. *Procedia CIRP*, 61, 335–340. <https://doi.org/10.1016/j.procir.2016.11.152>
160. van der Aalst, W. M. P. (2016). *Process mining: Data science in action* (2nd ed.). Springer. <https://doi.org/10.1007/978-3-662-49851-4>
161. van der Aalst, W. M. P., Adriansyah, A., & van Dongen, B. F. (2012). Replaying history on process models for conformance checking and performance analysis. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2(2), 182–192. <https://doi.org/10.1002/widm.1045>
162. Vasarhelyi, M. A., Alles, M. G., & Kogan, A. (2004). Principles of analytic monitoring for continuous assurance. *Journal of Emerging Technologies in Accounting*, 1(1), 1–21. <https://doi.org/10.2308/jeta.2004.1.1.1>
163. Verma, S., Dickerson, J., & Hines, K. (2024). Counterfactual explanations in machine learning: A survey. *ACM Computing Surveys*, 56(12), Article 1. <https://doi.org/10.1145/3677119>
164. Voss, M. D., & Williams, Z. (2013). Public-private partnerships and supply chain security: C-TPAT as a signal to regulators and partners. *Journal of Business Logistics*, 34(1), 1–12. <https://doi.org/10.1111/jbl.12030>
165. Wachter, S., Mittelstadt, B., & Floridi, L. (2017). Why a right to explanation of automated decision making does not exist in the GDPR. *International Data Privacy Law*, 7(2), 76–99. <https://doi.org/10.1093/idpl/ix005>
166. Waller, M. A., & Fawcett, S. E. (2013). Data science, predictive analytics, and big data: A revolution that will transform supply chain design and management. *Journal of Business Logistics*, 34(2), 77–84. <https://doi.org/10.1111/jbl.12010>
167. Wamba, S. F., Akter, S., Edwards, A., Chopin, G., & Gnanzou, D. (2015). How big data can make big impact: Findings from a systematic review. *International Journal of Production Economics*, 165, 234–246. <https://doi.org/10.1016/j.ijpe.2014.12.031>
168. Wang, L., Deng, T., Shen, Z. J. M., Hu, H., & Qi, Y. (2022). Digital twin driven smart supply chain. *Frontiers of Engineering Management*, 9(1), 56–70. <https://doi.org/10.1007/s42524-021-0186-9>
169. Wang, L., Wang, X. V., & Wang, Y. (2022). Digital twin driven smart supply chain. *Frontiers of Engineering Management*, 9(1), 56–70. <https://doi.org/10.1007/s42524-021-0186-9>
170. Wu, W., Zhao, Z., Shen, L., Kong, X. T. R., Guo, D., Zhong, R. Y., & Huang, G. Q. (2022). Just Trolley: Industrial IoT and digital twin enabled spatial temporal traceability and visibility for finished goods logistics. *Advanced Engineering Informatics*, 52, 101571. <https://doi.org/10.1016/j.aei.2022.101571>
171. Xia, L., Shanthikumar, J. G., & Zhu, S. (2020). Risk sensitive Markov decision processes with combined metrics: A mean variance approach. *Production and Operations Management*, 29(12), 2856–2876. <https://doi.org/10.1111/poms.13252>
172. Xu, L. D., Xu, E. L., & Li, L. (2018). Industry 4.0: State of the art and future trends. *International Journal of Production Research*, 56(8), 2941–2962. <https://doi.org/10.1080/00207543.2018.1444806>

173. Xu, L., Mak, S., Minaricova, M., & Brintrup, A. (2024). On implementing autonomous supply chains: A multi agent system approach. *Computers in Industry*, 161, 104120. <https://doi.org/10.1016/j.compind.2024.104120>
174. Yan, R., Sun, Z., & others. (2022). Reinforcement learning for transportation and logistics: A survey. *Transportation Research Part E: Logistics and Transportation Review*, 162, 102712. <https://doi.org/10.1016/j.tre.2022.102712>
175. Yang, Q., Liu, Y., Chen, T., & Tong, Y. (2019). Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology*, 10(2), Article 12. <https://doi.org/10.1145/3298981>
176. Yuan, X., He, P., Zhu, Q., & Li, X. (2019). Adversarial examples: Attacks and defenses for deep learning. *IEEE Transactions on Neural Networks and Learning Systems*, 30(9), 2805–2824. <https://doi.org/10.1109/TNNLS.2018.2886017>
177. Zhang, B., Tan, W. J., Cai, W., & Zhang, A. N. (2024). Leveraging multi agent reinforcement learning for digital transformation in supply chain inventory optimization. *Sustainability*, 16(22), 9996. <https://doi.org/10.3390/su16229996>
178. Zhang, K., Yang, Z., & Basar, T. (2021). Multi agent reinforcement learning: A selective overview of theories and algorithms. In *Handbook of Reinforcement Learning and Control* (pp. 321–384). https://doi.org/10.1007/978-3-030-60990-0_12
179. Zhang, Y., Chen, M., & Susilo, W. (2022). Information fusion for edge intelligence: A survey. *Information Fusion*, 78, 76–99. <https://doi.org/10.1016/j.inffus.2021.11.018>
180. Zhao, W., He, S., & Liu, C. (2023). State wise safe reinforcement learning: A survey. *Proceedings of the Thirty Second International Joint Conference on Artificial Intelligence (IJCAI 2023)*. <https://doi.org/10.24963/ijcai.2023/763>
181. Zhou, X., Chen, B., Gui, Y., & Cheng, L. (2025). Conformal prediction: A data perspective. *ACM Computing Surveys*, 57(1), Article 10. <https://doi.org/10.1145/3736575>
182. Zhou, Z. H., & Li, M. (2020). Tri training: Exploiting unlabeled data with robust evaluation. *IEEE Transactions on Knowledge and Data Engineering*, 32(5), 964–977. <https://doi.org/10.1109/TKDE.2019.2892626>