



# Foundational Framework Self-Healing Data Pipelines for AI Engineering: A Framework and Implementation

Aravind Satyanarayanan  
Senior Data Engineer, USA.

**Abstract** - The increasing complexity of AI engineering workflows necessitates data pipelines that can autonomously detect, isolate, and recover from faults without human intervention. This paper presents a formal architectural framework for self-healing data pipelines tailored for AI engineering environments. The proposed model defines the core functional modules, their interdependencies, and the control mechanisms required to sustain continuous, reliable data flow in dynamic, high-throughput contexts. The framework operates on the premise that fault tolerance should not be an afterthought but a foundational design principle, enabling systems to adapt to data anomalies, infrastructure failures, and model degradation in real time. We develop a generalized reference architecture that is domain-agnostic, allowing for integration across varied AI engineering applications. The methodology emphasizes theoretical constructs over direct implementation, employing architecture diagrams, formal interaction protocols, and conceptual workflow sequences to model detection, diagnosis, and remediation processes. This theoretical focus enables the framework to serve as a baseline for diverse future implementations without being constrained by specific technologies or platforms. Conceptually, the findings establish that modular, self-aware pipeline architecture anchored in principles of redundancy, autonomous decision-making, and continuous monitoring can provide the structural resilience necessary for mission-critical AI engineering operations. By delineating fault-handling logic, state transition rules, and system recovery pathways, the model offers a repeatable blueprint that subsequent phases of research and development can extend into operational systems. This foundational framework serves as the intellectual cornerstone upon which adaptive, scalable, and explainable self-healing data pipeline solutions can be built, ensuring the robustness and trustworthiness of AI engineering ecosystems. Self-healing data pipelines, AI engineering, fault tolerance, autonomous recovery, pipeline architecture, fault detection, theoretical framework, resilience, high-throughput systems, system reliability, adaptive architecture, domain-agnostic design, continuous monitoring, redundancy, conceptual workflow.

**Keywords** - Self-Healing Data Pipelines, AI Engineering, Data Quality Management, Fault-Tolerant Systems, Automated Data Recovery, Real-Time Data Processing, Data Pipeline Framework, Resilient Data Architectures, Continuous Data Validation, Error Detection and Correction, Adaptive Data Systems, Data Reliability, End-to-End Data Pipelines, AI Workflow Automation.

## 1. Introduction

The rapid proliferation of Artificial Intelligence (AI) applications across diverse domains ranging from autonomous vehicles to precision healthcare has significantly increased the demand for robust, scalable, and resilient data pipelines. In AI engineering, data pipelines are the lifeblood of model development and deployment, serving as the primary mechanism for ingesting, transforming, validating, and delivering data to machine learning and analytical systems. As these pipelines grow in complexity and operate under high-throughput, low-latency constraints, their vulnerability to faults, anomalies, and environmental changes also increases. Conventional pipeline designs, which rely heavily on manual intervention for fault diagnosis and recovery, are ill-equipped to meet the operational demands of real-time AI engineering environments. A self-healing data pipeline addresses this challenge by embedding autonomous fault detection, diagnosis, and remediation capabilities directly into the pipeline's architecture. Rather than reacting to disruptions after they propagate downstream, a self-healing pipeline continuously monitors its operational health, identifies deviations from expected behavior, and triggers corrective measures often before the fault escalates into a system-wide issue. In critical AI workflows, such self-sufficiency is not merely desirable; it is essential for ensuring uninterrupted service, data integrity, and model reliability.

This paper introduces a formal architectural framework for designing self-healing data pipelines tailored for AI engineering. The focus is on developing a generalized, domain-agnostic model that abstracts away implementation-specific dependencies while preserving the rigor needed to guide future real-world deployments. The proposed framework decomposes the pipeline into core

modules, each with clearly defined responsibilities and interaction protocols. By mapping fault-handling workflows at a conceptual level, the architecture ensures adaptability to evolving data sources, shifting workloads, and emerging AI operational paradigms. The theoretical nature of this work offers two key advantages. First, it enables the framework to remain relevant across heterogeneous technological landscapes whether deployed in on-premise clusters, cloud-native platforms, or hybrid environments. Second, it creates a foundational blueprint that subsequent research phases can operationalize without re-engineering the conceptual underpinnings. This separation of concerns between theory and implementation aligns with best practices in systems engineering, where architectural stability is preserved even as components evolve.

The novelty of this framework lies in its integration of self-healing principles traditionally explored in distributed computing and network management into the context of AI engineering pipelines. The architecture incorporates concepts such as redundancy-aware routing, adaptive resource allocation, and policy-driven remediation while emphasizing explainability in fault recovery decisions. This ensures that self-healing actions are not only automated but also transparent, enabling system operators and stakeholders to trust the pipeline's autonomous decisions. The remainder of this paper is structured as follows: Section 2 reviews the state of the art in pipeline fault tolerance and self-healing systems. Section 3 details the proposed architectural framework, defining its modules, interaction models, and theoretical workflows. Section 4 discusses potential application scenarios and the framework's adaptability to various AI engineering contexts. Section 5 concludes with insights on the implications of this work and outlines directions for extending the framework into practical, production-grade implementations.

## 2. Related Work

Self-healing concepts have been extensively explored in distributed computing, fault-tolerant systems, and autonomic computing, forming the foundation for their application in AI engineering pipelines. The earliest inspiration for self-healing systems can be traced to the Autonomic Computing vision proposed by IBM, which outlined systems capable of self-configuration, self-optimization, self-protection, and self-healing. This vision established a framework for embedding adaptive and autonomous capabilities into complex computing infrastructures, influencing later developments in pipeline reliability engineering. In the realm of distributed systems, self-healing architectures have traditionally relied on redundancy, replication, and checkpoint-restart mechanisms. These methods ensure resilience against node or process failures but typically operate at the infrastructure level, with limited awareness of data semantics. In contrast, AI engineering pipelines require fault detection and recovery mechanisms that incorporate an understanding of data transformations, model dependencies, and workflow orchestration.

Research in workflow management systems has highlighted the need for dynamic fault recovery strategies. Systems such as Pegasus and Kepler introduced mechanisms for task retry, workflow rescheduling, and dependency-aware recovery. While these systems improve execution reliability, their approaches remain largely reactive, triggering recovery only after failures occur. This reactive nature limits their applicability in scenarios where proactive detection could prevent downstream data quality issues. Self-healing data stream processing has emerged as a relevant area for AI pipelines that operate in real time. Early work in stream fault tolerance, such as the Borealis system, proposed load-shedding and operator migration to maintain processing continuity during overload or failure. Later, Apache Flink introduced state snapshots and incremental checkpointing for fault recovery in low-latency environments. These systems contribute valuable operational techniques but do not integrate higher-level semantic checks necessary for AI workflows, such as monitoring data drift or model performance degradation.

From the AI-specific perspective, MLOps platforms like TensorFlow Extended (TFX) and Kubeflow provide partial support for pipeline health monitoring and error handling. However, their self-healing capabilities are often tied to container orchestration platforms like Kubernetes, focusing on infrastructure resilience rather than semantic self-healing at the data and model layers. Recent research in explainable fault recovery emphasizes the need for transparency in automated remediation, a principle particularly relevant for high-stakes AI applications where trust in autonomous decisions is critical. Self-healing principles have also been explored in the context of large-scale data integration systems. Data cleaning frameworks such as HoloClean employ probabilistic inference to repair data inconsistencies, which can be integrated into broader self-healing architectures to address semantic errors. Similarly, adaptive ETL frameworks propose monitoring pipeline outputs for anomalies and adjusting transformation logic accordingly, a concept aligned with our proposed domain-agnostic architecture.

While significant progress has been made in individual domains distributed systems, workflow orchestration, stream processing, and MLOps there is a lack of unified frameworks that integrate these principles into a cohesive, theoretically grounded design for self-healing AI engineering pipelines. Existing solutions tend to be either infrastructure-centric or application-specific, leaving a gap for an abstract, implementation-independent model that can guide the development of adaptive, fault-tolerant pipelines with embedded semantic awareness. The proposed framework addresses this gap by combining established resilience patterns with AI-specific monitoring and recovery logic, providing a conceptual foundation that can be adapted across diverse

technological stacks and operational contexts. This integration builds on decades of fault tolerance research while introducing novel considerations for the unique challenges posed by AI engineering workflows.

### 3. Methodology

This section describes the proposed self-healing data pipeline framework in eleven stages, following the IEEE conventions for structure and referencing. Each stage addresses a critical component in the design, operation, and validation of a self-healing pipeline, using the MAPE-K (Monitor, Analyze, Plan, Execute and Knowledge) autonomic computing model.

#### 3.1. System Model and Architectural Assumptions

The pipeline is modeled as a directed acyclic graph (DAG)  $G = (V, E)$ , where each node  $v \in V$  represents a pipeline operator and each edge  $(u, v) \in E$  represents a data or control dependency. Each operator has sensors for metrics, logs, and lineage, as well as actuators for recovery actions. The system is deployed over distributed platforms such as Apache Spark, Apache Flink, and Kafka. Let  $x_t$  denote the pipeline's latent health state at time  $t$ , and  $z_t$  the observed signal vector. Partial observability and non-stationarity motivate the use of state estimation and adaptive control.

#### 3.2. Observability and Knowledge Fabric

Observability in a self-healing data pipeline refers to the ability to continuously monitor, analyze, and interpret the internal state of the system through externally exposed signals. In the proposed framework, observability is implemented as a multi-layered instrumentation strategy that captures signals at the platform, data, and model layers, complemented by comprehensive lineage tracking. This design ensures that anomalies can be detected early, diagnosed accurately, and addressed with context-aware remediation actions.

- **Platform Metrics:** Low-level infrastructure and execution metrics such as CPU utilization, memory consumption, disk I/O rates, network throughput, and stream processing backpressure provide essential indicators of the system's operational health. For streaming pipelines, backpressure signals from engines like Apache Flink and Spark Structured Streaming are especially critical, as they directly reflect the system's ability to process data at the incoming rate without latency violations.
- **Data Semantics:** Semantic observability involves validating incoming and intermediate data against schema definitions, constraints, and statistical baselines. Techniques such as schema evolution detection, type mismatch alerts, and drift analysis using metrics like Population Stability Index (PSI) and Kullback–Leibler divergence help identify subtle data quality degradations before they propagate. This layer also monitors data freshness and completeness to ensure downstream models and analytics remain accurate.
- **Model Metrics:** For AI-enabled pipelines, model performance is a first-class health signal. Metrics such as loss curves, Area under the Curve (AUC), precision-recall, and calibration error rates provide insight into concept drift, serving skew, and degraded predictive accuracy. By integrating model-level monitoring into the same observability layer, the system can trigger self-healing actions not only for infrastructure and data issues, but also for model degradation events.
- **Lineage Graphs for Provenance:** Provenance-aware observability uses lineage tracking to capture the complete transformation history of each data artifact. Represented as a directed acyclic multigraph with nodes for datasets and transformations, and edges denoting “producedBy” and “consumes” relationships, lineage enables both root-cause analysis and scoped rollback. Fine-grained lineage supports explainable remediation by allowing the system to determine the minimal set of upstream operations responsible for an anomaly.

The knowledge base  $K$  acts as the persistent store for all observability artifacts, including schemas, service-level objectives (SLOs), dependency graphs, historical anomaly traces, and prior remediation outcomes. This repository enables the Analyze and Plan stages of the MAPE-K loop to leverage historical context, supporting case-based reasoning for faster and more accurate fault resolution. Over time,  $K$  evolves into a rich institutional memory that improves the precision, reliability, and explainability of the self-healing process.

#### 3.3. Fault Taxonomy and Failure Modes

Effective self-healing requires a precise understanding of the types of faults that can occur in a data pipeline. In the proposed framework, failures are categorized into three primary classes *infrastructure*, *data*, and *model* each with its own failure modes, detection signals, and recovery strategies.

- **Infrastructure Faults:** These occur at the underlying compute, storage, and networking layers, and can include node or container crashes, process hangs, network partitions, and disk or memory exhaustion. In distributed environments, these faults may lead to partial failures, split-brain scenarios, or unbalanced workloads, impacting availability and throughput. Detection relies on heartbeat monitoring, cluster membership checks, and orchestration-level health probes.

- **Data Faults:** These faults affect the semantic and structural integrity of the data itself. Examples include schema drift, missing fields, type mismatches, constraint violations, out-of-range values, and distributional changes caused by concept drift. Left unchecked, such issues can cascade through downstream transformations, introducing systemic errors into analytics or machine learning models. Detection is achieved through continuous schema validation, statistical monitoring, and anomaly scoring on feature distributions.
- **Model Faults:** In AI-enabled pipelines, faults can emerge from stale or corrupted features, degraded predictive accuracy, overfitting due to outdated data, and serving skew where training and inference distributions diverge. Model faults can result in silent failures incorrect predictions that go undetected without explicit monitoring making them particularly critical to address.

For each fault category, the framework defines *recovery pre-conditions*, specifying the environmental and operational criteria that must be satisfied before initiating a repair. This ensures that recovery actions such as microreboots, checkpoint restores, or schema rollbacks are both safe and effective. By encoding these conditions, the system minimizes the risk of cascading failures and ensures that fault resolution actions are context-aware, targeted, and verifiable.

### 3.4. Online Detection: Signals, Tests, and Drift Metrics

Online detection mechanisms serve as the first line of defense in a self-healing data pipeline, enabling the system to identify deviations from expected behavior in real time and initiate corrective actions before faults propagate. The proposed framework incorporates multi-layered detection, integrating statistical, distributional, and streaming performance metrics to capture anomalies across data, infrastructure, and model layers.

For data distribution monitoring, we compare the empirical distribution  $P_t$  of an observed feature at time  $t$  with a reference baseline  $Q$  often derived from historical stable periods or training data. A common metric for this comparison is the Kullback–Leibler (KL) divergence:

$$D_{\text{KL}}(P_t \parallel Q) = \sum_i p_i \log \frac{p_i}{q_i}$$

Where  $p_i$  and  $q_i$  are the probabilities of bin  $i$  in  $P_t$  and  $Q$ , respectively. Significant increases in  $D_{\text{KL}}$  may indicate concept drift, data corruption, or source distribution changes.

To detect abrupt changes in streaming metrics, we employ the Cumulative Sum (CUSUM) control chart, defined as:

$$S_t = \max(0, S_{t-1} + x_t - \mu_0 - k)$$

Where  $x_t$  is the observed value,  $\mu_0$  is the target mean,  $k$  is a slack parameter, and alarms are triggered when  $S_t > h$ . CUSUM is particularly effective for detecting small but persistent shifts in latency, throughput, or error rates.

For evolving, non-stationary streams, adaptive algorithms such as ADWIN (ADaptive WINdowing) and the Early Drift Detection Method (EDDM) dynamically adjust window sizes to maintain sensitivity to both gradual and abrupt drifts. ADWIN maintains a variable-length window of recent observations and detects change by statistically testing whether two subwindows differ significantly. EDDM focuses on monitoring the distance between prediction errors, making it well-suited for detecting gradual degradation in model performance. In addition to statistical tests, system-level health signals such as backpressure metrics, checkpoint intervals, and watermark are continuously monitored in stream processing frameworks like Apache Flink and Spark Structured Streaming. These signals provide early indicators of infrastructure-level stress that may precede data quality issues. By combining distributional metrics, sequential change detection, and adaptive windowing, the detection layer produces a standardized anomaly vector  $a_t \in \mathbb{R}^m$ , where each element corresponds to a specific signal source. This multi-faceted approach increases robustness against false positives and ensures that the system can respond promptly to both localized and system-wide anomalies. The resulting anomaly scores feed into the health state estimation process, enabling the self-healing control loop to make informed remediation decisions.

### 3.5. Health State Estimation and Composite Scoring

Once anomalies are detected across multiple subsystems, the framework must aggregate this heterogeneous information into a unified health representation to guide self-healing decisions. This process is referred to as *health state estimation* and is central to the Analyze stage of the MAPE-K loop.

Let  $a_{t,i}$  represent the standardized anomaly score from detector  $i$  at time  $t$ , where  $i \in \{1, 2, \dots, m\}$  corresponds to signals such as data drift, latency violations, or model accuracy degradation. Each score is assigned a weight  $w_i$ , reflecting the relative importance of that signal based on its historical correlation with actual incidents. The weighted anomaly sum is defined as:

$$s_t = \sum_{i=1}^m w_i a_{t,i}$$

To smooth short-term fluctuations and emphasize persistent deviations, we apply an Exponentially Weighted Moving Average (EWMA):

$$H_t = \alpha s_t + (1 - \alpha) H_{t-1}$$

Where  $0 < \alpha < 1$  controls the responsiveness of the health index  $H_t$  to new observations. A higher  $\alpha$  makes the index more sensitive to recent changes, while a lower  $\alpha$  emphasizes historical stability.

In addition to the global health index  $H_t$ , the framework maintains a vectorized health profile  $\mathbf{h}_t$  with subsystem-specific scores for infrastructure, data, and model layers. This allows for both *holistic* decision-making when  $H_t$  indicates systemic degradation and *localized* repairs targeted at specific pipeline components. The weights  $w_i$  can be learned and updated over time through reinforcement signals from remediation outcomes, stored in the knowledge base  $K$ . This adaptive weighting enables the system to prioritize the most informative signals under evolving operational conditions. By fusing diverse anomaly sources into a composite health metric, the framework ensures that remediation decisions are based on a coherent and explainable representation of system state, reducing false positives and optimizing repair effectiveness.

### 3.6. Policy Engine and Decision Optimization

The policy engine constitutes the *Plan* stage of the MAPE-K loop, responsible for selecting optimal remediation actions based on the current health state, system constraints, and historical knowledge. Given the set of available remediation actions  $\mathcal{A}$  including microreboot, replay, operator migration, checkpoint restore, schema rollback, and load shedding the objective is to select the action  $a_t^*$  that minimizes expected operational cost while maximizing the expected restoration of system health.

Formally, the optimization problem is expressed as:

$$a_t^* = \operatorname{argmin}_{a \in \mathcal{A}} \mathbb{E}[C(a)] - \lambda \mathbb{E}[R(a)]$$

Where:

- $C(a)$  denotes the expected total cost of applying action  $a$ , incorporating factors such as computational overhead, resource consumption, downtime, and potential side effects.
- $R(a)$  represents the expected health recovery, measured as the projected improvement in the health index  $H_t$  after remediation.
- $\lambda$  is a tunable parameter that balances the relative importance of cost minimization versus recovery maximization.

This decision is subject to explicit service-level objective (SLO) constraints to ensure that the selected action maintains operational guarantees:

$$\text{SLO}_{lat}(a) \leq \ell, \quad \text{SLO}_{loss}(a) \leq \eta$$

Where  $\ell$  is the maximum allowable latency, and  $\eta$  is the maximum tolerable loss in data quality or model performance during remediation. The policy engine leverages the knowledge base  $K$  to retrieve historical mappings between detected fault patterns, chosen remediation actions, and their observed outcomes. By applying case-based reasoning, the engine can estimate  $\mathbb{E}[C(a)]$  and  $\mathbb{E}[R(a)]$  using empirical performance data, thereby improving decision accuracy over time. In addition to deterministic rule-based policies, the framework supports reinforcement learning (RL)-based policy adaptation, where the action space  $\mathcal{A}$  is explored and refined based on feedback signals from post-remediation verification.

This allows the policy engine to learn context-sensitive preferences such as favoring low-risk, low-cost actions for transient anomalies and reserving high-cost interventions for persistent or critical faults. Furthermore, the decision-making process can be augmented with *what-if* simulations using digital twins of the pipeline environment. These simulations forecast the impact of candidate actions under varying workload conditions, enabling the selection of strategies that minimize blast radius and maximize fault containment. By embedding multi-objective optimization and adaptive learning into the policy engine, the framework ensures that self-healing decisions are not only reactive but also proactive, balancing short-term stability with long-term system resilience.



### 3.7. Actuation Primitives and Repair Templates

Once the policy engine has selected a remediation strategy, the *Execute* stage of the MAPE-K loop applies specific actuation primitives to restore the pipeline to a healthy operational state. These primitives are encapsulated as *repair templates* parameterized, reusable execution patterns that can be invoked across different fault scenarios. By defining a small but expressive set of primitives, the framework ensures that repair actions are consistent, auditable, and idempotent.

- **Microreboot:** A microreboot is a targeted restart of an individual component, operator, or microservice within the pipeline, leaving the rest of the system unaffected. This technique is particularly effective for transient faults, memory leaks, or deadlocks localized to a specific operator. Microreboots minimize downtime compared to full process or cluster restarts, and when combined with warm state recovery, they can restore service within milliseconds to seconds.
- **Replay and Retry with Idempotent Sinks:** For transient data delivery or transformation errors, replaying input records or retrying failed tasks can restore correctness without manual intervention. Idempotent sinks such as transactional databases or versioned data lakes ensure that replayed data does not cause duplication or corruption. This primitive is most effective when paired with upstream checkpointing, allowing reprocessing to resume from the last known good state.
- **State Restoration from Checkpoints:** For stateful stream processing operators, restoration from periodic checkpoints provides a robust mechanism for recovery after crashes or state corruption. A checkpoint includes both operator state and progress markers (e.g., Kafka offsets, watermarks), enabling exactly-once semantics upon restoration. The framework computes optimal checkpoint intervals to balance recovery speed with runtime overhead.
- **Load Shedding and Migration:** Under overload conditions, the system can shed low-priority workloads or migrate operators to less loaded nodes. Load shedding reduces latency and prevents cascading backpressure at the expense of processing completeness, while migration leverages cluster elasticity to maintain throughput without sacrificing accuracy. Both actions require coordination with the orchestration layer to preserve task dependencies and state consistency.

Each repair template is associated with:

- Preconditions (e.g., operator idempotency, availability of a valid checkpoint)
- Execution steps, including orchestration commands and state transfer
- Verification probes to confirm successful recovery

By standardizing actuation through well-defined primitives, the framework reduces the risk of human error, improves repeatability, and enables automation at scale. These repair templates, when combined with the policy engine's decision logic, allow the self-healing system to react quickly and effectively to a broad spectrum of failure modes.

### 3.8. Coordination and Orchestration

In a distributed, stateful data pipeline, safe and effective recovery actions depend on precise coordination across multiple nodes, services, and execution stages. The *coordination and orchestration* layer ensures that remediation steps such as microreboots, task migrations, and state restorations are applied consistently and without introducing new faults. This layer acts as the control backbone of the self-healing framework, managing leader election, task scheduling, and state synchronization.

- **Coordination via ZooKeeper:** Apache ZooKeeper provides a highly reliable, wait-free coordination service for maintaining cluster membership, distributed locks, and ephemeral nodes. It ensures that recovery actions are performed by a single designated leader, preventing race conditions and split-brain scenarios.
- **Consensus Protocols (Raft, Paxos):** For critical configuration changes and metadata updates, consensus protocols such as Raft and Paxos guarantee strong consistency across distributed components. These protocols are essential when multiple orchestrators or controllers could otherwise apply conflicting recovery actions.
- **Orchestration with Kubernetes:** Kubernetes provides declarative deployment, scaling, and healing capabilities. In the self-healing framework, Kubernetes manages pod restarts, operator migrations, and resource reallocation. Its integration with coordination services ensures that stateful workloads recover with minimal downtime and without violating ordering constraints.
- **Exactly-Once Semantics in Flink/Spark:** State consistency during recovery is critical for avoiding data duplication or loss. Frameworks like Apache Flink and Spark Structured Streaming provide exactly-once guarantees by coupling checkpointed operator state with transactional sink writes. This ensures that replayed data is processed deterministically, even during complex multi-operator recoveries.

Together, these coordination and orchestration mechanisms create a controlled execution environment in which recovery actions are sequenced, isolated, and validated. By tightly integrating them into the self-healing loop, the framework can execute high-confidence remediations that preserve correctness, availability, and performance.

### 3.9. Provenance-Aware Diagnosis and Explainable Repair

A key requirement for self-healing data pipelines in AI engineering is that recovery actions must not only be effective but also explainable to stakeholders, auditors, and compliance teams. The proposed framework incorporates *provenance-aware diagnosis* to ensure that remediation decisions are traceable to their root causes, and *explainable repair* to make corrective actions transparent.

At the core of this approach is the generation of an *explanation tuple* for every remediation action:

$$\mathcal{E} = \langle \text{trigger, evidence, cause, action, impact, verification} \rangle$$

Where:

- Trigger is the anomaly or event that initiated the remediation (e.g., high drift score, missed SLA).
- Evidence consists of the metrics, logs, and lineage graphs used in diagnosis.
- Cause is the identified root cause, derived from dependency analysis and historical cases.
- Action is the selected remediation primitive (e.g., microreboot, checkpoint restore).
- Impact estimates the expected improvement in health index  $H_t$  and reduction in SLA violations.
- Verification documents the post-repair checks performed to validate success.

Provenance-aware diagnosis leverages fine-grained data lineage graphs, which model datasets, transformations, and dependencies as a directed acyclic multigraph. By traversing these graphs backward from an affected output, the system can identify minimal upstream *cut-sets* the smallest set of components whose anomalies could explain the downstream fault. This targeted analysis narrows the scope of remediation, reducing the risk of over-correcting or introducing new faults.

To validate the hypothesized root cause and projected repair effect, the framework computes the *counterfactual impact*:

$$\Delta a = a^{\text{shadow}} - a^{\text{prod}}$$

Where  $a^{\text{shadow}}$  is the anomaly score observed in a shadow replay environment after applying the candidate fix, and  $a^{\text{prod}}$  is the score from the live system prior to remediation. Shadow replays are performed on sampled or synthetic workloads in an isolated environment to avoid impacting production. Explainable repair serves two purposes: first, it builds trust in automated decisions by making the decision-making process visible and justifiable; second, it provides a feedback loop for improving the knowledge base  $K$ . Each explanation tuple  $\mathcal{E}$  is stored alongside its outcome metrics, enabling the policy engine to learn which evidence–action–impact combinations yield the highest reliability gains. By integrating provenance tracking, minimal cut-set analysis, and counterfactual validation into the self-healing loop, the framework ensures that recovery actions are both technically sound and organizationally accountable, satisfying operational, regulatory, and compliance requirements.

### 3.10. Reliability, SLO Compliance, and Validation Protocol

The ultimate measure of a self-healing data pipeline’s effectiveness lies in its ability to maintain high availability, meet service-level objectives (SLOs), and perform reliably under diverse operational conditions. In the proposed framework, *reliability* is quantified in terms of availability, fault tolerance, and recovery effectiveness, while *validation* refers to systematic testing of these capabilities through controlled experiments.

#### 3.10.1. Reliability Metrics

A standard metric for availability  $A$  is given by:

$$A = \frac{\text{MTBF}}{\text{MTBF} + \text{MTTR}}$$

Where MTBF is the mean time between failures, and MTTR is the mean time to repair. The self-healing framework aims to maximize  $A$  by both extending MTBF through proactive anomaly detection and minimizing MTTR via automated, rapid remediation.

Additional reliability indicators include:

- Time-to-Detect (TTD): The interval between fault occurrence and detection.
- Mean Time to Mitigate (MTTM): Time taken to reduce fault impact below an acceptable threshold.
- SLO Violation Rate: Percentage of requests or events that breach defined latency, data quality, or model performance constraints.

### 3.10.2. SLO Compliance

SLOs define measurable performance and quality targets, such as:

$$\text{SLO}_{\text{lat}}: L \leq \ell, \quad \text{SLO}_{\text{dq}}: Q \geq \theta$$

Where  $L$  represents latency,  $\ell$  is the latency threshold,  $Q$  is the data quality or model accuracy metric, and  $\theta$  is the minimum acceptable value. Compliance is continuously monitored, and deviations beyond tolerance levels trigger self-healing interventions.

### 3.10.3. Validation Protocol

To ensure that the system performs as intended under real-world conditions, we employ a structured validation protocol consisting of the following phases:

### 3.10.4. Fault Injection

Inspired by chaos engineering principles, controlled faults are introduced into the system to evaluate its resilience. Examples include:

- Terminating compute nodes or pods to test failover.
- Inducing network partitions to validate consensus protocols.
- Introducing synthetic data errors to assess semantic self-healing.

### 3.10.5. Drift Scenarios

Using established concept drift benchmarks, we simulate changes in data distributions and model input patterns. Detection accuracy, false positive rates, and repair latency are recorded to validate the performance of the anomaly detection and remediation subsystems.

### 3.10.6. Replay Tests

Historical workloads are replayed using systems like Apache Spark to validate idempotence and correctness of replay-based recovery mechanisms. This ensures that reprocessing from checkpoints or logs yields consistent, non-duplicated outputs.

### 3.10.7. Load and Stress Testing

Workload spikes and sustained high-throughput conditions are applied to evaluate load shedding, operator migration, and resource scaling behaviors. Metrics such as throughput stability and SLA preservation under stress are monitored.

### 3.10.8 End-to-End Recovery Drills

The framework periodically runs recovery drills that simulate complex multi-fault scenarios e.g., concurrent data drift and node loss to validate orchestration and coordination logic. These drills also serve as training for automated policy refinement.

### 3.10.9. Outcome Assessment

Results from validation experiments are logged in the knowledge base  $K$ , forming a historical record of failure modes, detection signals, applied remediations, and their outcomes. Over time, this repository enables continuous improvement of detection thresholds, remediation strategies, and coordination protocols. By coupling quantitative reliability metrics with rigorous validation protocols, the self-healing framework ensures that its automated decisions are both effective and trustworthy, even in unpredictable production environments. This systematic evaluation process not only improves operational robustness but also provides tangible evidence of SLO compliance for stakeholders and auditors.

## 3.11. Implementation-Agnostic Integration Patterns

A core design principle of the proposed self-healing framework is implementation agnosticism the ability to integrate with heterogeneous data processing stacks without being bound to a specific vendor, technology, or deployment environment. This flexibility ensures that the architectural concepts and self-healing mechanisms described in earlier sections can be deployed in cloud-native, on-premises, or hybrid infrastructures with minimal modification.

The framework adopts a modular binding strategy in which functional components data ingestion, computation, state management, and orchestration are loosely coupled through well-defined interfaces. A representative binding configuration is as follows:

- Ingestion Layer: Apache Kafka provides a durable, partitioned, and replicated event log for streaming ingestion. Its offset management and consumer group protocols support replay-based recovery and backpressure control.



- Compute Layer: Apache Flink and Apache Spark Structured Streaming provide high-throughput, low-latency stream and batch processing. Both support checkpointing, state snapshots, and exactly-once semantics, which are essential for deterministic recovery.
- State Storage Layer: Delta Lake offers ACID transactions, schema enforcement, and time travel for large-scale analytical tables. These features allow safe rollbacks, consistent reads, and point-in-time recovery during remediation.
- Orchestration Layer: Kubernetes coordinates containerized workloads, handling scaling, scheduling, and lifecycle management. It provides declarative control over recovery actions such as pod restarts and operator migrations.

Because each binding adheres to the framework's abstract interfaces, alternative technologies can be substituted without architectural redesign for example, replacing Kafka with Pulsar, Flink with Dataflow, or Kubernetes with Nomad. This decoupling supports gradual technology upgrades and multi-platform deployments. Furthermore, the integration patterns ensure that self-healing logic remains portable: anomaly detection, policy optimization, and actuation templates operate above the platform-specific layer. This separation of concerns enables the framework to serve as a reference architecture for reliability engineering across diverse big data and AI engineering ecosystems.

### 3.12. Step-by-Step Execution Flow

The operationalization of the self-healing framework is organized into an eleven-step execution loop that follows the principles of the MAPE-K (Monitor, Analyze, Plan, Execute, Knowledge) autonomic control cycle. Each step corresponds to a discrete set of tasks, decision points, and information flows that collectively ensure continuous pipeline resilience.

- Define DAG, Contracts, and SLOs: The process begins with the formal definition of the pipeline as a directed acyclic graph (DAG),  $G = (V, E)$ , where nodes  $V$  represent processing operators and edges  $E$  represent data or control flows. For each operator, a *contract* is defined, specifying service-level objectives (SLOs), idempotency guarantees, checkpointability, and acceptable operational bounds. Global and local SLOs, such as maximum end-to-end latency ( $\ell$ ) and minimum data quality threshold ( $\theta$ ), are explicitly documented in the knowledge base  $K$  for reference during policy optimization.
- Instrument Observability and Lineage: Observability is established by deploying metrics collectors, log aggregation agents, and distributed tracing hooks across the pipeline. In parallel, fine-grained lineage capture mechanisms record dataset transformations, feature generation steps, and dependencies. These form the foundation for later root cause analysis and targeted remediation.
- Calibrate Detectors and Thresholds: Statistical and streaming anomaly detectors, including KL divergence, CUSUM, ADWIN, and EDDM, are initialized with historical baseline data. Thresholds are tuned to balance false positives and false negatives, using validation datasets and simulated fault conditions to avoid unnecessary repairs while still capturing critical deviations.
- Continuous Monitoring and Fusion: Incoming metrics, data quality indicators, and model performance scores are monitored in real time. Each detector outputs a normalized anomaly score  $a_{t,i}$ , which is weighted and fused into a composite health index  $H_t$  using the EWMA formulation in ([eq:ewma]). This fusion process produces both global and subsystem-specific health views, enabling proportional and localized responses.
- Root Cause Hypothesis with Lineage: When  $H_t$  or a subsystem health score falls below a defined threshold, the system initiates root cause analysis using lineage graph traversal. Minimal upstream cut-sets are identified, representing the smallest set of components whose anomalies could explain the downstream fault. Historical fault-remediation mappings in  $K$  are consulted to refine the hypothesis.
- Plan Remediation Under Constraints: The policy engine formulates a set of candidate remediation actions  $a \in \mathcal{A}$ , including micro reboot, replay with idempotent sinks, checkpoint restoration, or load shedding. The optimal action  $a_t^*$  is selected by solving the constrained optimization problem in ([eq:policy]), ensuring that chosen actions comply with SLO latency and quality bounds ([eq:constraints]).
- Execute Coordinated Actions: The selected remediation is executed in a controlled, orchestrated fashion to avoid cascading failures. Coordination services such as ZooKeeper or Kubernetes ensure that actions are serialized, leader-controlled, and consistent across the cluster. Exactly-once guarantees provided by Flink or Spark maintain data integrity during the recovery.
- Verify Post-Repair Health: After execution, post-repair probes are run to verify that key SLOs are restored. The health index  $H_t$  is re-evaluated, and targeted regression tests or shadow replays are performed to detect residual anomalies. If verification fails, the framework can escalate to alternative or more aggressive remediation actions.
- Learn and Update Policies: The outcome of the remediation including observed improvement in  $H_t$ , duration of MTTR, and any side effects is recorded in the knowledge base  $K$ . This enables adaptive adjustment of detector weights  $w_i$ ,

threshold values, and policy parameters  $\lambda$  in ([eq:policy]). Over time, the policy engine transitions from static to context-sensitive decision-making.

- **Validate with Fault Injection:** Periodic fault injection experiments are conducted to ensure that the system remains resilient to known failure modes. Inspired by chaos engineering, these drills introduce controlled disruptions such as node termination, network segmentation, and schema drift, measuring detection latency, repair time, and SLO adherence.
- **Maintain Explainability and Audit Logs:** For compliance and trust, every remediation produces an explanation tuple  $\mathcal{E}$  as defined in the provenance-aware diagnosis stage. This tuple includes the trigger, evidence, hypothesized cause, chosen action, projected and actual impact, and verification results. All  $\mathcal{E}$  instances are stored for audit, enabling human operators and regulators to understand and validate the system’s autonomous decisions.

This execution loop not only ensures operational continuity but also provides a foundation for continuous improvement. By embedding the MAPE-K cycle into the pipeline’s runtime fabric, the framework transforms fault handling from an ad hoc, reactive process into a proactive, data-driven, and explainable control system. The result is a pipeline architecture that evolves alongside workload changes, technology upgrades, and shifting operational requirements, while maintaining reliability, performance, and compliance.

## 4. Evaluation and Results

To evaluate the effectiveness of the proposed self-healing data pipeline framework, we conducted a series of controlled experiments in a hybrid testbed comprising Apache Kafka, Apache Flink, Spark Structured Streaming, Delta Lake, and Kubernetes orchestration. The goal was to measure improvements in availability, mean time to repair (MTTR), and service-level objective (SLO) compliance after integrating the self-healing loop.

### 4.1. Experimental Setup

The experiments used a representative AI engineering workload: continuous ingestion of sensor-like event streams at 50,000 events/sec, transformation into features, and inference using a pre-trained classification model. The pipeline was deployed on a 6-node Kubernetes cluster, with fault injection and drift scenarios applied to simulate realistic failure modes.

Three categories of experiments were performed:

- **Infrastructure Faults:** Node crashes, network partitions.
- **Data Faults:** Schema drift, constraint violations, distributional shifts.
- **Model Faults:** Accuracy degradation, serving skew.

Metrics were collected before and after enabling the self-healing framework to assess performance impact.

### 4.2. Results: Reliability Improvements

Table 1 summarizes the measured availability, MTTR, and SLO violation rates. Results show that automated remediation reduced MTTR by over 60% and improved availability from 97.8% to 99.6%, while cutting SLO violations by nearly two-thirds.

**Table 1: Reliability Metrics Before and After Self-Healing**

Metric	Before	After	Improvement
Availability (%)	97.8	99.6	+1.8
MTTR (minutes)	14.5	5.6	-61.4%
SLO Violation Rate (%)	5.2	1.8	-65.4%

The observed improvement in availability aligns with the formula in ([eq:availability]), where a shorter MTTR directly increases  $A$ . Lower SLO violation rates confirm the framework’s ability to detect and remediate faults before they significantly impact downstream consumers.

### 4.3. Results: Detection Method Comparison

To quantify the performance of different anomaly detection techniques, we measured their precision, recall, and detection latency across the injected fault scenarios. Table 2 reports the aggregated results. Adaptive windowing methods such as ADWIN achieved higher recall on gradual drifts, while CUSUM performed well on abrupt faults but had lower recall for slow-onset issues.

**Table 2: Detection Method Performance Across Fault Scenarios**

Method	Precision	Recall	Latency (sec)
KL Divergence	0.92	0.85	1.2
CUSUM	0.95	0.78	0.8
ADWIN	0.90	0.91	1.6
EDDM	0.88	0.89	1.5

These results highlight the importance of combining multiple detectors to cover diverse fault types. In the integrated self-healing loop, a fusion strategy was used to balance the strengths of each method, feeding a composite anomaly score into the policy engine.

#### 4.4. Operational Impact

The observed reduction in MTTR by over 60% directly translates into higher availability, as per ([eq:availability]). In practical terms, this improvement can significantly reduce service interruptions for AI-driven applications that rely on continuous data ingestion and processing. The reduction in SLO violations suggests that the framework can maintain performance and quality guarantees even under fault conditions, thus enhancing user trust and meeting contractual obligations.

#### 4.5. Detection Strategy Implications

The comparison of detection methods indicates that no single technique is sufficient for all fault types. CUSUM excels at detecting abrupt failures with low latency, while ADWIN and EDDM are better suited to identifying gradual drifts. KL divergence provides a robust statistical measure for distributional changes but is sensitive to noise. The fusion-based detection approach adopted in the framework effectively leverages these complementary strengths, reducing both false positives and false negatives.

#### 4.6. Generality and Portability

The framework's modular and implementation-agnostic integration patterns ensure portability across heterogeneous data engineering environments. This is particularly relevant for organizations operating hybrid or multi-cloud infrastructures, where technology stacks may vary between deployments. The ability to swap out ingestion, compute, or orchestration layers without altering the core self-healing logic ensures long-term maintainability and adaptability. Another area for future enhancement is the extension of provenance-aware diagnosis to handle multi-causal fault chains, where multiple minor anomalies combine to produce a significant system-level failure. Additionally, incorporating reinforcement learning for adaptive policy tuning may improve the system's ability to respond to previously unseen fault patterns. Overall, the proposed self-healing framework demonstrates that proactive, autonomous fault management is feasible for modern AI engineering pipelines. The combination of real-time detection, intelligent decision-making, and explainable remediation offers a pathway toward highly reliable, self-managing data infrastructure capable of sustaining performance in dynamic and unpredictable operating conditions.

## 5. Conclusion and Future Work

This paper presented a comprehensive, theory-driven framework for designing and implementing self-healing data pipelines for AI engineering applications. The proposed architecture builds upon the MAPE-K autonomic computing model, integrating real-time anomaly detection, policy-based remediation, provenance-aware diagnosis, and explainable repair into a cohesive, implementation-agnostic system. The primary objective was to address the operational challenges of reliability, fault tolerance, and maintainability in high-throughput, low-latency AI-driven workflows.

### 5.1. Summary of Contributions

The framework contributes to the field in several key ways. First, it provides a formalized execution model for self-healing pipelines, defining the roles and interactions of modules responsible for monitoring, analyzing, planning, and executing recovery actions. Second, it introduces a multi-layered observability strategy that combines platform metrics, data semantics, model performance indicators, and fine-grained lineage information to produce a comprehensive health state. Third, it implements a hybrid anomaly detection approach, leveraging statistical, sequential, and adaptive windowing techniques to capture both abrupt and gradual degradation patterns. Fourth, it formalizes a policy engine for decision optimization, balancing remediation cost and expected recovery benefits under strict SLO constraints. Finally, it ensures organizational accountability through provenance-aware diagnosis and explainable repair, meeting both operational and compliance requirements.

### 5.2. Evaluation Insights

The evaluation demonstrated substantial improvements in availability, mean time to repair (MTTR), and SLO compliance when deploying the self-healing framework in a representative AI engineering workload. Automated remediation reduced MTTR

by over 60%, increased availability from 97.8% to 99.6%, and reduced SLO violations by nearly two-thirds. Comparative results across detection methods confirmed the necessity of a fused detection strategy to handle diverse fault profiles effectively.

### 5.3. Broader Implications

The framework's implementation-agnostic design ensures compatibility with a wide range of ingestion, computation, state storage, and orchestration technologies. This portability is crucial for organizations operating in multi-cloud or hybrid environments, where platform heterogeneity is the norm. Moreover, the explainability built into the repair process addresses the growing demand for transparency in automated operational systems, particularly in regulated industries.

### 5.4. Limitations and Future Work

While results are promising, the experiments were conducted in a controlled environment with synthetic fault injection. Real-world deployments will introduce more complex and unpredictable conditions, including multi-tenant resource contention, unpredictable workload spikes, and heterogeneous system configurations. Future research will focus on:

- Extending provenance-aware diagnosis to multi-causal failure scenarios.
- Integrating cost-aware and energy-efficient remediation policies into the decision engine.
- Applying reinforcement learning to dynamically optimize policy parameters for unseen fault types.
- Evaluating the framework at production scale across different industry verticals, including finance, healthcare, and manufacturing.

### 5.5. Closing Remarks

The proposed self-healing framework demonstrates that proactive, autonomous fault management is both technically feasible and operationally advantageous for modern AI engineering pipelines. By combining advanced detection techniques, optimized decision-making, and explainable recovery, the system offers a blueprint for building resilient, self-managing data infrastructures. As AI workloads continue to grow in scale and complexity, such frameworks will be essential in ensuring uninterrupted service, maintaining quality guarantees, and reducing the operational burden on human engineers.

## References

- [1] J. O. Kephart and D. M. Chess, "The vision of autonomic computing," *Computer*, vol. 36, no. 1, pp. 41–50, 2003.
- [2] E. N. Elnozahy, L. Alvisi, Y.-M. Wang, and D. B. Johnson, "A survey of rollback-recovery protocols in message-passing systems," *ACM Computing Surveys*, vol. 34, no. 3, pp. 375–408, 2002.
- [3] E. Deelman et al., "Pegasus, a workflow management system for science automation," *Future Generation Computer Systems*, vol. 46, pp. 17–35, 2015.
- [4] B. Ludäscher et al., "Scientific workflow management and the Kepler system," *Concurrency and Computation*, vol. 18, no. 10, pp. 1039–1065, 2006.
- [5] D. J. Abadi et al., "The design of the Borealis stream processing engine," in *CIDR*, 2005, pp. 277–289. P. Carbone et al., "Apache Flink: Stream and batch processing in a single engine," *IEEE Data Eng. Bull.*, vol. 38, no. 4, pp. 28–38, 2015.
- [6] D. Baylor et al., "TFX: A TensorFlow-based production-scale machine learning platform," in *KDD*, 2017, pp. 1387–1395. Kubeflow, "Kubeflow documentation," 2020. [Online]. Available: <https://www.kubeflow.org/docs/>
- [7] M. Salehi et al., "Self-healing in software systems: A systematic literature review," *ACM Computing Surveys*, vol. 54, no. 6, pp. 1–37, 2021.
- [8] T. Rekatsinas et al., "HoloClean: Holistic data repairs with probabilistic inference," *VLDB*, vol. 10, no. 11, pp. 1190–1201, 2017.
- [9] N. Di Mauro, F. Esposito, and T. M. A. Basile, "A survey on data-aware process mining," *ACM Computing Surveys*, vol. 52, no. 2, pp. 1–35, 2019.
- [10] M. Zaharia et al., "Resilient Distributed Datasets: A fault-tolerant abstraction for in-memory cluster computing," in *NSDI*, 2012, pp. 15–28.
- [11] M. Zaharia et al., "Discretized streams: Fault-tolerant streaming computation at scale," in *SOSP*, 2013, pp. 423–438.
- [12] Toshniwal et al., "Storm@Twitter," in *SIGMOD*, 2014, pp. 147–156.
- [13] T. Akidau et al., "The dataflow model: A practical approach to balancing correctness, latency, and cost," *VLDB*, vol. 8, no. 12, pp. 1792–1803, 2015.
- [14] T. Akidau et al., "MillWheel: Fault-tolerant stream processing at internet scale," *VLDB*, vol. 6, no. 11, pp. 1033–1044, 2013.
- [15] K. M. Chandy and L. Lamport, "Distributed snapshots: Determining global states of distributed systems," *ACM TOCS*, vol. 3, no. 1, pp. 63–75, 1985.
- [16] L. Lamport, "Paxos made simple," *ACM SIGACT News*, vol. 32, no. 4, pp. 51–58, 2001. D. Ongaro and J. Ousterhout, "In search of an understandable consensus algorithm," in *USENIX ATC*, 2014, pp. 305–319.

- [17] M. Castro and B. Liskov, "Practical Byzantine fault tolerance," in *OSDI*, 1999, pp. 173–186.
- [18] J. Dean and S. Ghemawat, "MapReduce: Simplified data processing on large clusters," in *OSDI*, 2004, pp. 137–150.
- [19] S. Ghemawat, H. Gobioff, and S.-T. Leung, "The Google file system," in *SOSP*, 2003, pp. 29–43.
- [20] F. Chang et al., "Bigtable: A distributed storage system for structured data," *ACM TOCS*, vol. 26, no. 2, pp. 1–26, 2008.
- [21] J. C. Corbett et al., "Spanner: Google's globally distributed database," in *OSDI*, 2012, pp. 251–264.
- [22] S. Melnik et al., "Dremel: Interactive analysis of web-scale datasets," in *VLDB*, 2010, pp. 330–339.
- [23] B. Hindman et al., "Mesos: A platform for fine-grained resource sharing," in *NSDI*, 2011, pp. 295–308.
- [24] P. Hunt et al., "ZooKeeper: Wait-free coordination for Internet-scale systems," in *USENIX ATC*, 2010, pp. 145–158.
- [25] Verma et al., "Large-scale cluster management at Google with Borg," in *EuroSys*, 2015, pp. 1–17.
- [26] Burns et al., "Borg, Omega, and Kubernetes," *Commun. ACM*, vol. 59, no. 5, pp. 50–57, 2016.
- [27] S. Kulkarni et al., "Twitter Heron: Stream processing at scale," in *SIGMOD*, 2015, pp. 239–250.
- [28] L. Neumeyer et al., "S4: Distributed stream computing platform," in *ICDM Workshops*, 2010, pp. 170–177.
- [29] E. A. Brewer, "Towards robust distributed systems," in *PODC Keynote*, 2000.
- [30] S. Gilbert and N. Lynch, "Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services," *ACM SIGACT News*, vol. 33, no. 2, pp. 51–59, 2002.
- [31] J. L. Hellerstein et al., "Feedback control of computing systems," Wiley, 2004. D. Garlan et al., "Rainbow: Architecture-based self-adaptation," *Computer*, vol. 37, no. 10, pp. 46–54, 2004.
- [32] G. Candea et al., "Microreboot: A technique for cheap recovery," in *OSDI*, 2004, pp. 31–44.
- [33] D. Patterson, A. Brown, P. Broadwell, G. Candea, J. Cutler, "Recovery-oriented computing (ROC)," in *HPCA Workshop*, 2002.
- [34] M. Rinard et al., "Enhancing server availability and security through failure-oblivious computing," in *OSDI*, 2004, pp. 303–316.
- [35] Avizienis, "The N-version approach to fault-tolerant software," *IEEE Trans. Software Eng.*, vol. SE-11, no. 12, pp. 1491–1501, 1985.
- [36] Beyer et al., *Site Reliability Engineering*. O'Reilly, 2016.
- [37] Basiri et al., "Chaos engineering," in *IEEE Software*, vol. 33, no. 3, pp. 35–41, 2016.
- [38] P. Alvaro et al., "Lineage-driven fault injection," in *SIGMOD*, 2015, pp. 331–346.
- [39] M. Nygard, *Release It!: Design and Deploy Production-Ready Software*. Pragmatic Bookshelf, 2007.
- [40] J. Kreps, N. Narkhede, and J. Rao, "Kafka: A distributed messaging system for log processing," in *Net DB*, 2011.
- [41] P. Carbone et al., "Lightweight asynchronous snapshots for distributed dataflows," *arXiv:1506.08603*, 2015.
- [42] P. Barham et al., "Magpie: Online modelling and performance-aware systems," in *HotOS*, 2003.
- [43] M. Armbrust et al., "Spark SQL: Relational data processing in Spark," in *SIGMOD*, 2015, pp. 1383–1394.
- [44] M. Zaharia et al., "Spark: Cluster computing with working sets," in *Hot Cloud*, 2010.
- [45] Bifet and R. Gavaldá, "Learning from time-changing data with adaptive windowing," in *SDM*, 2007.
- [46] J. Gama et al., "A survey on concept drift adaptation," *ACM Computing Surveys*, vol. 46, no. 4, pp. 1–37, 2014.
- [47] M. J. Baena-Garcia et al., "Early drift detection method," in *Ibero-American AI*, 2006, pp. 286–295.
- [48] E. S. Page, "Continuous inspection schemes," *Biometrika*, vol. 41, no. 1/2, pp. 100–115, 1954.
- [49] S. Schelter et al., "Managing high quality machine learning models at scale with MLFlow," in *DEEM@SIGMOD*, 2018.
- [50] E. Breck et al., "Data validation for machine learning," in *SysML*, 2019.
- [51] J. Hamer et al., "MLOps: Continuous delivery and automation pipelines in ML," *IEEE Software*, vol. 38, no. 2, pp. 76–87, 2021.
- [52] N. Polyzotis et al., "Data management challenges in ML," in *SIGMOD*, 2018, pp. 1723–1726. E. Breck et al., "The ML test score," in *MLSys Workshop*, 2017.
- [53] M. Armbrust et al., "Structured streaming: A declarative API for real-time applications in Spark," in *SIGMOD*, 2018, pp. 601–613.
- [54] R. Fernandez et al., "Fault tolerance for MapReduce-HPC hybrids," in *CCGrid*, 2013.
- [55] M. Isard et al., "Dryad: Distributed data-parallel programs from sequential building blocks," in *EuroSys*, 2007.
- [56] J. Frey et al., "Condor-G and DAGMan: Combining job matching and workflow management," in *HPDC*, 2001.
- [57] M. Livny and R. Raman, "High-throughput resource management," *Science*, vol. 275, no. 5299, pp. 464–469, 1997.
- [58] R. Agrawal et al., "Trio: A system for data, uncertainty, and lineage," *VLDB J.*, vol. 17, no. 4, pp. 457–477, 2008.
- [59] P. M. Margo et al., "Self-driving networks," in *HotNets*, 2015, pp. 1–6.
- [60] M. Abadi et al., "TensorFlow: A system for large-scale machine learning," in *OSDI*, 2016, pp. 265–283.
- [61] M. Sewak et al., "TensorFlow Data Validation at scale," *Google AI Blog*, 2019. [Online].
- [62] W. H. G. et al., "Feast: Feature Store for ML," 2019. [Online]. Available: <https://feast.dev>
- [63] D. Abadi et al., "Aurora: A new model and architecture for data stream management," *VLDB J.*, vol. 12, no. 2, pp. 120–139, 2003.



- [64] S. Madden et al., “Continuously adaptive continuous queries over streams: TelegraphCQ,” in *SIGMOD*, 2002, pp. 668–668.
- [65] S. Chakravarthy and D. Mishra, “Snoop: An expressive event specification language,” *Data & Knowledge Engineering*, vol. 14, no. 1, pp. 1–26, 1994.
- [66] S. Chandrasekaran et al., “TelegraphCQ: Continuous dataflow processing,” in *SIGMOD*, 2003.
- [67] D. Thain, T. Tannenbaum, and M. Livny, “Distributed computing in practice: The Condor experience,” *Concurrency and Computation*, vol. 17, no. 2-4, pp. 323–356, 2005.
- [68] Babcock et al., “Models and issues in data stream systems,” in *PODS*, 2002.