



Optimizing Data Quality in Real-Time: A Self-Healing Pipeline Approach

Aravind Satyanarayanan
Senior Data Engineer, USA.

Abstract: The reliability of AI-driven decision-making systems depends not only on the robustness of their architectures but also on the consistent quality of the data they process. In real-time analytics environments, ensuring high data quality is often in tension with meeting stringent latency requirements. This paper introduces a theoretical framework for optimizing data quality in self-healing data pipelines by employing a quantitative decision model that balances latency constraints with adaptive data validation. The approach is grounded in formalizing the optimization problem through explicit constraints on processing time, detection rates, and acceptable error margins. At the core of this model is an adaptive validation function capable of dynamically tuning its verification intensity based on observed data distributions and system performance metrics. Rather than relying on fixed, rule-based checks that may either underperform during data drift or overburden the system under high load, the proposed method continuously calibrates itself to achieve optimal trade-offs. Using simulated data generated from synthetic probability distributions, we evaluate the model's behavior under varying levels of noise, drift, and system stress.

Our findings indicate that adaptive validation strategies consistently outperform static validation rules in non-stationary environments, enabling pipelines to maintain high data fidelity without compromising throughput. The theoretical results also identify threshold conditions under which adaptive checks provide the greatest benefit, offering a decision-making guide for system architects. By embedding this optimization model within a self-healing pipeline, we enhance not only its ability to detect and repair anomalies but also to proactively sustain the quality of streaming data in mission-critical applications. This work contributes to the growing body of theory that positions data quality assurance as an integral, quantitative component of resilient AI infrastructure, with broad applicability across finance, healthcare, cybersecurity, and other latency-sensitive sectors. Self-healing data pipelines, real-time data quality, adaptive validation, latency optimization, quantitative decision model, streaming analytics, anomaly detection, data drift, non-stationary environments, AI infrastructure.

Keywords: Data Quality, Real-Time Data Processing, Self-Healing Pipelines, Data Pipeline Optimization, Automated Data Validation, Streaming Data, Data Governance, Fault-Tolerant Systems, Data Monitoring, Data Reliability, Error Detection and Correction, Data Engineering, Adaptive Data Pipelines.

1. Introduction

The increasing reliance on artificial intelligence (AI) and machine learning (ML) in critical decision-making has elevated the role of data pipelines from auxiliary infrastructure to a core determinant of system reliability. In domains such as financial markets, healthcare diagnostics, cybersecurity monitoring, and autonomous systems, the timeliness and accuracy of data directly influence the validity of downstream analytics and predictions. As such, modern data engineering faces a dual imperative: to maintain high data quality while satisfying stringent real-time processing constraints. This challenge is magnified in high-velocity streaming environments where data distributions are non-stationary, sources are heterogeneous, and workloads are subject to sudden surges. Traditional approaches to data validation often based on static rules and fixed thresholds are insufficient in these dynamic contexts. While static validation can be computationally efficient, it risks either failing to detect subtle but significant anomalies or introducing unnecessary processing delays when the rules are overly conservative. Moreover, in systems where the cost of delayed data delivery is high, any additional latency can erode the operational value of the data. This necessitates a shift from static, rule-bound validation toward adaptive, context-aware strategies that can adjust validation intensity and methodology in real time.

Self-healing pipelines, designed to detect, diagnose, and remediate faults autonomously, offer a promising foundation for this shift. However, while considerable research has addressed the architectural resilience of such pipelines focusing on failure detection, automated recovery, and fault-tolerant orchestration comparatively less attention has been devoted to embedding formal, quantitative models for data quality optimization within them. Without a principled mechanism to balance validation rigor against latency budgets, even the most resilient pipelines can deliver results that are timely but untrustworthy, or accurate but too late to act upon. This paper introduces a mathematical decision model for real-time data quality optimization within self-healing pipelines. The model formulates the trade-off between latency and quality as a constrained optimization problem, incorporating parameters such as maximum allowable processing delay, anomaly detection rate, and acceptable error tolerance. Central to the approach is an

adaptive validation function that adjusts its operation dynamically based on current system conditions and statistical properties of the incoming data stream. The function leverages observed distributions to calibrate the intensity of validation checks, ensuring that computational resources are allocated efficiently and that validation overhead scales proportionally to data risk.

To evaluate the theoretical viability of the model, we conduct simulation experiments using synthetic probability distributions to represent varying data characteristics and operational conditions. These simulations explore system behavior under scenarios involving data drift, noise injection, and variable load intensity, allowing us to identify threshold conditions where adaptive validation offers maximal benefit over static rule sets. The findings demonstrate that adaptive validation can consistently maintain high data fidelity while respecting latency constraints, particularly in environments subject to frequent distributional changes.

The contributions of this work are threefold:

- A formal decision-theoretic framework for optimizing real-time data quality in self-healing pipelines.
- An adaptive validation function capable of scaling its rigor to match operational risk in streaming contexts.
- A theoretical performance analysis via simulation, offering system architects actionable insights into parameter tuning for latency-quality trade-offs.

By integrating quantitative data quality optimization into the design of self-healing pipelines, this research advances the broader goal of creating AI-ready infrastructure capable of delivering both timeliness and trustworthiness in mission-critical applications.

2. Related Work

Our framework builds on three strands of prior work: (i) performance and control foundations for streaming systems, (ii) data management and transactional robustness, and (iii) data quality, anomaly detection, and MLOps.

- Performance and control. Classic capacity and queueing results inform our latency budgets and headroom policies. Amdahl's law and the Universal Scalability Law model parallel efficiency and contention effects, while queueing theory provides principled latency-throughput trade-offs under bursty arrivals. Feedback control and elastic provisioning underpin our MAPE-K loop and risk-to-budget mapping. At the execution layer, data-parallel and graph-parallel engines motivate our operator design and partition-aware validation (Dryad, Naiad/Timely, Differential Dataflow, PowerGraph, GraphX). Streaming analytics and seasonal/robust statistics (STL, Hampel influence functions, Anscombe's quartet) guide our drift and plausibility proxies.
- Data management and transactions. Our fail-closed, auditable remediation aligns with decades of work on recovery and coordination. Sagas, ARIES logging, and two-phase commit provide templates for idempotent replay, partial rollbacks, and consistency under failures. Coordination avoidance highlights when global synchronization can be safely reduced a principle we exploit via hard/soft partitioning of checks. In storage and warehousing, relational foundations, OLAP design, B-tree engineering, and benchmark-driven evaluation (TPC-DS) inform our lineage indexing and table-format choices; modern log/columnar substrates (Kafka, Spark SQL, Delta Lake) shape our implementation targets.
- Data quality, anomaly detection, and MLOps. Foundational work on assessing and filtering data quality motivates our multi-dimensional quality vector and policy guardrails. Robust detection across logs and streams leverages surveys and benchmarks (NAB) and modern detectors (DeepLog, MIDAS) that inspire our proxy and RCA libraries. On the ML lifecycle side, technical-debt analyses, continuous integration for models, production-readiness rubrics, feature synthesis, and data-management challenges frame our learning-based utility estimation and audit requirements. Finally, the antifragility perspective underscores our use of shadow-mode fault injection to improve policies under perturbations.

In summary, our contribution operationalizes these lines into a unified, risk-sensitive, self-healing control layer that allocates validation effort under latency constraints, explains and audits its actions, and continually learns utility/cost surfaces on top of widely deployed streaming substrates.

3. Methodology

This section presents a structured methodology for *Optimizing Data Quality in Real-Time: A Self-Healing Pipeline Approach*. The proposed method is anchored in a closed-loop control paradigm that integrates real-time monitoring, quantitative decision-making, and automated remediation, enabling continuous optimization of data quality without compromising latency-sensitive workloads. The approach is designed to be implementation-agnostic, ensuring applicability across a wide range of modern streaming architectures. The core of the methodology lies in a quantitative decision model that dynamically balances the rigor of data validation against stringent end-to-end latency constraints. Incoming data streams are continuously profiled to extract key quality metrics, such as completeness, accuracy, consistency, and timeliness. These metrics are compared against predefined

service-level objectives (SLOs) to assess compliance in real time. When deviations are detected, a risk assessment module evaluates the potential impact on downstream consumers, prioritizing issues based on severity and business context.

Following the assessment, the planning component selects appropriate adaptive validation and remediation strategies. These strategies may include schema evolution handling, missing value imputation, outlier correction, or source re-ingestion, all executed within bounded latency budgets. The execution layer applies these corrective actions in a controlled manner, with formal correctness guarantees where applicable, while maintaining transactional consistency. The methodology incorporates a feedback subsystem that records all detected anomalies, applied remediations, and resulting quality improvements. This historical feedback is used to refine validation thresholds, improve anomaly detection accuracy, and update remediation policies over time, effectively creating a self-learning system. By decoupling the methodology from specific technologies, the design can be mapped to common streaming stacks involving ingestion (e.g., Kafka), compute (e.g., Flink, Spark Structured Streaming), state storage (e.g., Delta Lake), orchestration (e.g., Kubernetes), and monitoring frameworks, ensuring extensibility and operational resilience.

3.1. System Model and Problem Statement

We consider a high-velocity streaming data pipeline tasked with processing an ordered sequence of records $\{x_t\}$, each arriving at time t with an associated metadata vector m_t . The metadata vector encapsulates essential contextual information such as the originating source identifier, schema version, data lineage pointer, processing watermark, and the observed arrival delay. This information is critical for assessing both operational and semantic aspects of the stream in real time. The pipeline is constrained by a strict latency budget ℓ_{\max} that applies to each record or to aggregated micro-batch windows W_t , while simultaneously ensuring that the data quality delivered to downstream consumers does not fall below a predefined target threshold q_{\min} .

To formalize quality assessment, we define a multi-dimensional quality vector for each processing window W_t :

$$\mathbf{Q}_t = (Q_t^{\text{valid}}, Q_t^{\text{complete}}, Q_t^{\text{consist}}, Q_t^{\text{fresh}}, Q_t^{\text{accuracy}}),$$

Where each component captures a distinct quality dimension. In practice, these dimensions may rely on proxy metrics when ground truth labels are unavailable. Examples of such proxies include constraint satisfaction ratios for validity, referential integrity rates for completeness, drift-adjusted plausibility scores for accuracy, and prediction consistency agreements for semantic stability. The system dynamically estimates these proxies in real time to infer quality without halting the streaming process. We denote by Φ the set of n available validation operators (or checks), indexed by $i \in \{1, \dots, n\}$. For each window W_t , the pipeline can select an execution intensity $v_{i,t} \in [0,1]$ for operator i , representing a fractional application of its full computational and validation capability. Each operator i is characterized by a latency cost function $c_i(v_{i,t})$ and an *expected quality gain* $\Delta q_i(v_{i,t} | \mathcal{S}_t)$, conditioned on the state summary \mathcal{S}_t of the stream. The state summary incorporates compact statistical representations, distributional drift indicators, lineage trust scores, and recent violation histories, thus serving as the decision context for adaptive validation.

The objective of the pipeline at each window W_t is to select the set of validation intensities $\{v_{i,t}\}$ that maximizes the net utility U_t defined as the total expected quality gain minus a latency penalty subject to an instantaneous processing budget:

$$\begin{aligned} \max_{\{v_{i,t}\}} \quad & U_t = \sum_{i=1}^n \Delta q_i(v_{i,t} | \mathcal{S}_t) - \lambda_t \sum_{i=1}^n c_i(v_{i,t}) \\ \text{s.t.} \quad & \sum_{i=1}^n c_i(v_{i,t}) \leq B_t, \quad 0 \leq v_{i,t} \leq 1. \end{aligned}$$

Here, λ_t is an adaptive Lagrange multiplier or penalty weight that dynamically tunes the trade-off between validation thoroughness and processing latency, reflecting system-wide priorities at time t . The instantaneous budget B_t is not static; rather, it is *risk-driven*, computed by a self-healing controller that continuously integrates system load, observed anomaly rates, and downstream criticality to adjust the permissible validation workload in real time. This formulation captures the essential challenge of *adaptive data quality optimization in streaming systems*: balancing computational resource allocation against dynamic quality risks under stringent latency constraints. By structuring the problem in this manner, the framework provides a mathematically rigorous foundation for implementing self-healing pipelines capable of sustaining both performance and trustworthiness in volatile data environments.

3.2. Quality Dimensions and Proxies

To operationalize the proposed framework in a streaming environment, the data quality components are instantiated using streaming-amenable estimators that allow for constant-time updates and bounded memory usage. Each dimension is associated with a set of well-defined checks and a corresponding proxy metric that can be computed online.

- **Validity:** This dimension captures whether incoming records adhere to the expected schema and data type specifications, as well as to semantic constraints such as permissible ranges, enumerated categorical vocabularies, and pattern conformity (e.g., regular expressions for identifier formats). The proxy for validity is defined as the fraction of records within the current processing window that satisfy all active constraints. This metric is updated incrementally as each record is ingested.
- **Completeness:** Completeness assesses the presence of all required fields and the overall coverage of expected data items within a given window. It accounts for non-nullness, as well as window-level coverage against expected counts derived from metadata or prior historical patterns. In streaming contexts, delayed arrivals are compensated for using watermarking techniques. The proxy metric is a coverage ratio, adjusted by quantiles of watermark delay to avoid penalizing temporarily late but valid records.
- **Consistency:** Consistency refers to the degree to which records conform to cross-entity constraints, such as referential integrity between related tables, functional dependencies, and conditional dependencies within and across streams. It also includes temporal consistency checks between correlated data sources. The proxy is the proportion of satisfied dependency tests over the set of tests applicable to the records in the current window.
- **Freshness:** Freshness measures the timeliness of records, specifically the inter-arrival delays and the lag between event time and processing time. Stale data can degrade analytical accuracy and operational decision-making. The proxy for freshness is computed as the moving quantiles of $(t_{\text{process}} - t_{\text{event}})$, complemented by staleness flags when delays exceed configured thresholds.
- **Accuracy:** In many real-time scenarios, ground truth labels are unavailable at the moment of data arrival. Accuracy is therefore approximated through plausibility checks, invariance tests, prediction agreement among redundant or ensemble models, and statistical distance measures comparing the incoming data distribution to trusted reference profiles. The proxy metric combines drift-adjusted likelihood scores with ensemble agreement rates to quantify probable correctness.

To maintain computational efficiency, all proxies are updated using *streaming sketches* and *incremental estimators*. For example, Count–Min sketches efficiently approximate categorical frequency distributions, exponential moving averages track rate-based metrics, and online quantile summaries estimate delay distributions without retaining full data histories. This approach ensures that quality estimation scales gracefully with throughput while preserving responsiveness under strict latency budgets.

3.3. Self-Healing Control Loop

The self-healing controller is designed following the classical *MAPE-K* (Monitor–Analyze–Plan–Execute with Knowledge) autonomic control paradigm, adapted to the constraints and requirements of real-time streaming data quality optimization. The loop operates continuously at the granularity of a processing window or micro-batch, dynamically adjusting validation intensities and remediation actions based on observed conditions.

- **Monitor:** The monitoring phase is responsible for collecting low-cost yet high-value telemetry from all relevant components of the streaming pipeline. This includes per-operator latencies, input and output queue depths, watermark lags, and violation counters for each quality dimension. Additionally, the system computes quality proxy metrics (as defined in the preceding section) and collects statistical drift signals that indicate distributional shifts. The monitoring subsystem employs lightweight, streaming-friendly data structures to avoid adding significant overhead, ensuring that observation does not materially impact latency budgets.
- **Analyze:** In the analysis phase, the collected telemetry is aggregated into a concise state summary \mathcal{S}_t . From this, the controller computes a *risk score* r_t that reflects both the probability and the potential impact of quality degradation if the current state persists without intervention. This risk score is derived from statistical models or rule-based heuristics that map observed anomalies (e.g., increased schema violations, rising drift scores) to estimated downstream harm. The computation of r_t may integrate probabilistic reasoning, Bayesian updating, or learned regression models trained on historical violation–impact pairs. The result is a single scalar or multi-dimensional risk vector that serves as the primary decision driver for subsequent planning.
- **Plan:** The planning phase determines how to allocate the available validation budget B_t and select operator intensities $\{v_{i,t}\}$ to maximize the utility function defined in Equation [eq:budget]. The budget B_t itself is risk-sensitive: in high-risk situations, the controller may allocate more time and resources to intensive validation, whereas in low-risk states it can relax validation to preserve throughput. The optimization process may be solved via exact methods (e.g., constrained quadratic programming) or via heuristics for faster adaptation. In addition to budget allocation, the planning stage

constructs a remediation plan for violations already detected. This plan is governed by policy constraints, such as compliance requirements, allowable data transformations, and the cost–benefit trade-offs of quarantining versus repairing affected records.

- **Execute:** In the execution phase, the controller applies the chosen validation intensities by dynamically tuning operator configurations (e.g., sampling rates, rule activation levels). It triggers remediation strategies, which may include quarantining suspect records for offline inspection, applying automated repair functions (such as imputing missing values or correcting type mismatches), or rerouting data to alternative processing paths. Routing decisions may be used to isolate problematic data sources or to switch to redundant streams when available. Execution is followed by the logging of both actions and outcomes into the knowledge base to ensure traceability and auditability.

The Knowledge component serves as the long-term memory of the controller, storing mappings from contextual states to effective actions based on past experience. It maintains calibration parameters for utility estimation, allowing the risk–budget trade-off to be refined over time. Additionally, it preserves lineage-aware explanations of past remediation actions, enabling the system to provide human operators with auditable justifications for its decisions. This lineage-awareness is crucial for compliance in regulated domains, where data modifications must be accompanied by verifiable reasoning chains. By integrating continuous monitoring, context-aware analysis, adaptive planning, and auditable execution within the MAPE-K framework, the self-healing control loop enables the pipeline to maintain high-quality data delivery in dynamic, high-throughput environments. This architecture supports both proactive interventions anticipating quality degradation before it impacts consumers and reactive remediation when violations occur, thereby ensuring resilience under fluctuating workloads and evolving data characteristics.

3.4. Risk-Driven Budget Allocation

A central element of the proposed self-healing pipeline is the dynamic allocation of the validation budget B_t based on the estimated risk level r_t . This mechanism ensures that validation resources are spent proportionally to the likelihood and severity of potential quality degradation, thereby balancing throughput preservation and quality assurance.

The mapping from risk to budget is defined as:

$$r_t = \sigma(\mathbf{w}^\top \mathbf{z}_t), \quad B_t = B_{\min} + \kappa r_t \Delta B_{\max},$$

Where $\sigma(\cdot)$ denotes the logistic sigmoid function, producing a normalized risk score in $[0,1]$. The feature vector \mathbf{z}_t aggregates a diverse set of context indicators that collectively characterize the health of the incoming data stream and the operational environment. Representative features include:

- **Distribution Drift Distances:** Statistical distances (e.g., Jensen–Shannon divergence, Wasserstein distance) between recent data windows and long-term reference profiles, capturing gradual or abrupt shifts in feature distributions.
- **Anomaly Rates:** The fraction of records violating active validation rules or exceeding anomaly thresholds, normalized over the observation window.
- **Source Trust Priors:** Historical reliability scores for each data source, derived from past violation frequencies and correction accuracy.
- **Schema Change Indicators:** Flags for detected changes in schema versions or field-level data types, which often precede elevated violation rates.
- **SLO Slack:** The margin between current end-to-end latency and the maximum allowed ℓ_{\max} , indicating the available headroom for additional validation.

The parameter B_{\min} enforces a baseline level of validation, guaranteeing that a minimal set of essential checks is always executed, even under conditions of very low estimated risk. The term ΔB_{\max} represents the maximum additional budget that can be allocated when risk is at its peak, while κ acts as a sensitivity coefficient controlling how aggressively the budget scales with r_t . This allocation policy provides two key benefits. First, it prevents over-validation in low-risk scenarios, preserving system throughput and avoiding unnecessary latency. Second, it ensures that during high-risk periods such as after a schema migration, a detected drift spike, or a sudden rise in anomaly rates the system can immediately expand validation coverage and intensity without manual intervention. By tying B_t directly to measurable and explainable features, the controller maintains transparency and auditability in its decision-making, satisfying both operational efficiency and compliance requirements.

3.5. Adaptive Validation Operators

In the proposed self-healing pipeline, validation is implemented through a family of operators, each capable of adjusting its execution parameters via an *intensity* control variable $v \in [0,1]$. This design enables fine-grained tuning of computational effort

versus quality gain, allowing the system to dynamically scale the rigor of checks based on the current budget B_t and assessed risk r_t .

- **Syntactic and Type Checks:** These operators verify schema conformance, type integrity, and pattern compliance. Intensity controls include the sampling rate of records inspected, the per-field depth of type-checking (for example, validating nested JSON structures fully versus partially), and early termination upon detecting initial violations (short-circuiting). At low intensity, only a small sample of records is inspected; at high intensity, every record undergoes deep validation.
- **Constraint Validation:** Enforces domain-specific constraints such as functional dependencies, conditional dependencies, and referential integrity rules. Intensity can be varied by activating only a subset of constraints, reducing the dependency-mining window size, or adjusting the frequency of probabilistic constraint tests. For example, under latency pressure, the system might validate only the most violation-prone constraints, while under low load, it can activate the full constraint set.
- **Statistical Profiling:** Monitors distributions of numeric and categorical features to detect drift or anomalies. Adjustable parameters include histogram bin counts, sketch width for streaming frequency estimation, update frequency, and significance thresholds for statistical drift tests (e.g., Kolmogorov–Smirnov or Jensen–Shannon). Lower intensities use coarser summaries and infrequent updates, while higher intensities maintain finer-grained statistics with stricter detection thresholds.
- **Cross-Stream Consistency:** Validates alignment between related streams, such as matching transaction records with corresponding audit logs. Intensity controls include join sampling rates, minimum key coverage targets, and tolerances for temporal alignment discrepancies. At low intensity, only a subset of keys is cross-checked; at high intensity, the operator enforces near-complete coverage and tighter alignment bounds.
- **Semantic Validation:** Incorporates model-in-the-loop plausibility checks, where a predictive model assesses whether record values are contextually reasonable. Intensity adjustments include the proportion of records passed to the model, the complexity of the model invoked (e.g., lightweight classifier versus deep neural network), and thresholds for flagging anomalies. This ensures that expensive model calls are used sparingly under tight budgets but extensively when quality risk is high.

For each operator i , we either precompute or learn a cost model $c_i(v)$ representing expected computational overhead, and a marginal utility function $\Delta q_i(v)$ capturing the improvement in quality metrics from increasing intensity. When historical performance data is available, these models can be parameterized offline using regression or simulation. In cases where explicit cost–utility models are unavailable, we adopt *multi-armed bandit* strategies to estimate $\Delta q_i(v)$ online. This involves controlled exploration of different intensities over time to identify the most cost-effective operating points under varying contexts \mathcal{S}_t . This adaptive framework allows the pipeline to continuously optimize its validation strategy in response to evolving data characteristics, operational constraints, and risk conditions, ensuring both resource efficiency and sustained data quality in real-time environments.

3.6. Optimization and Selection

The allocation of validation intensities $\{v_{i,t}\}$ for each operator family is driven by the constrained optimization problem in ([eq:budget]). In practice, solving this problem exactly at every processing window W_t may be computationally prohibitive in high-throughput settings. Therefore, we adopt an incremental, near-optimal heuristic based on the *benefit-to-cost* ratio, defined for each operator i at a given intensity v as:

$$\rho_i(v) = \frac{\partial \Delta q_i(v)}{\partial v} / \frac{\partial c_i(v)}{\partial v}.$$

Here, $\Delta q_i(v)$ represents the marginal improvement in quality metrics obtained by increasing intensity v , while $c_i(v)$ denotes the corresponding latency or computational cost.

The greedy algorithm proceeds by identifying the operator with the maximum $\rho_i(v)$ and incrementally increasing its intensity $v_{i,t}$ in small quanta. This process continues until the aggregate cost reaches the allocated budget B_t , subject to policy-enforced minima for mandatory checks. Such constraints ensure that essential validations (e.g., schema compliance or PII redaction) are never disabled, even under severe budget pressure. When $\Delta q_i(v)$ exhibits diminishing returns and $c_i(v)$ is convex, this greedy allocation closely approximates the true optimum. For scenarios involving interdependent operators where the benefit of one operator depends on the activation of another we adopt a two-stage selection procedure. In the first stage, the algorithm identifies a feasible base set that satisfies all hard constraints and dependency requirements. In the second stage, the residual budget is distributed among *soft-utility* operators to maximize total utility.

To prevent instability in operator intensities caused by fluctuating real-time risk scores, we introduce a *hysteresis* mechanism:

$$v_{i,t} \leftarrow \beta v_{i,t-1} + (1 - \beta) \hat{v}_{i,t}, \quad \beta \in [0,1),$$

where $\hat{v}_{i,t}$ is the raw output of the optimizer at time t . The smoothing factor β controls the inertia of adjustments, dampening sudden changes in intensity. Additionally, a cooldown timer enforces a minimum period between significant parameter shifts, reducing oscillations when risk scores hover near decision thresholds. This combined approach greedy selection, constraint enforcement, interdependency handling, and stabilization enables the controller to make fast, explainable allocation decisions that remain robust under dynamic load and risk conditions. By balancing responsiveness with stability, the pipeline maintains both high data quality and predictable latency performance in real-time operation.

3.7. Streaming State Estimation

The state summary \mathcal{S}_t serves as a compact, continuously updated representation of the pipeline's operational context. It enables both risk modeling and utility prediction for adaptive validation, without imposing prohibitive computational or storage overheads. Maintaining \mathcal{S}_t in real-time requires streaming-compatible algorithms that can operate under strict latency and memory constraints. The state is composed of the following components:

- **Distribution Profiles:** The statistical characteristics of incoming records are maintained through reservoir or stratified sampling, ensuring representative coverage even under skewed arrival patterns. Frequency distributions and quantile estimates are captured using streaming sketches such as Count–Min Sketch for categorical variables and GK (Greenwald–Khanna) summaries for numerical data. This allows rapid detection of unusual shifts in value distributions while avoiding full-data scans.
- **Drift Signals:** Concept and data drift are monitored via rolling two-sample tests between a reference window and the most recent data window. To support real-time processing, we employ streaming-friendly divergence measures such as Jensen–Shannon divergence or Population Stability Index (PSI) computed incrementally. Adaptive window sizing is applied so that drift detection remains sensitive under both stable and volatile regimes, with dynamic thresholds informed by historical variance.
- **Lineage Trust:** Each data source is assigned a reliability score that evolves over time through Bayesian credit assignment. Observed constraint violations, late arrivals, and confirmed incident reports reduce the trust score, while periods of clean, timely data increase it. This metric acts as a prior in the risk model, allowing higher scrutiny of historically unreliable sources without penalizing well-behaved streams.
- **SLO Slack:** The available performance headroom is estimated from end-to-end latency measurements and real-time queue dynamics. By applying Little's Law ($L = \lambda W$) to arrival and service rates, the system calculates the slack margin before breaching the latency budget ℓ_{\max} . This slack estimate enables the controller to decide whether additional validation checks can be executed without exceeding service-level objectives.

These streaming estimators are designed to be composable, lightweight, and implementation-agnostic, allowing deployment across heterogeneous streaming architectures. The resulting \mathcal{S}_t is fed into both the risk-driven budget allocation policy and the per-operator utility estimation models, ensuring that validation intensity decisions are informed by an accurate, continuously updated operational picture.

4. Violation Detection and Root Cause Analysis

When observed violations exceed the pre-defined tolerance thresholds, the analyzer component of the self-healing control loop initiates a structured root cause analysis (RCA) workflow. This workflow is designed to operate under real-time constraints while maintaining a high degree of diagnostic accuracy. The process is divided into three key stages: localization, correlation, and causal hypothesis generation each of which is supported by streaming-compatible analytics and knowledge-base augmentation.

- **Localization:** The first step in RCA is to pinpoint the exact location and nature of the violation. The system identifies failing constraints (e.g., schema mismatch, referential integrity failures, or business rule violations) by comparing observed proxy metrics against expected thresholds. Field-level granularity is used to isolate the specific attributes exhibiting anomalous behavior. Lineage pointers are then leveraged to trace the offending data segments back through the pipeline to their originating sources and processing stages. This enables not only the isolation of problematic datasets but also the identification of transformations or joins that may have introduced the defect.
- **Correlation:** Once the violation is localized, the next stage is to investigate its temporal and contextual associations with other events or signals. The analyzer examines time-series patterns to detect spikes in violations and aligns them with potential triggers such as recent schema deployments, version upgrades in upstream services, data ingestion delays, or changes in load distribution. The correlation process incorporates diverse telemetry streams, including drift metrics, anomaly detection alerts, upstream error logs, and infrastructure performance counters. By integrating multi-modal

signals, the analyzer reduces the likelihood of false attributions and builds a richer context for the eventual hypothesis generation.

- **Causal Hypotheses:** The final stage involves formulating candidate explanations for the observed violations. The analyzer uses both rule-based reasoning and probabilistic inference to generate a ranked list of hypotheses, each with an associated confidence score. Examples of such hypotheses include: (a) schema evolution in the upstream system without a corresponding update in the downstream parser, (b) delayed partition arrival caused by backpressure in the message broker, (c) corruption of records due to an upstream sampling policy change, or (d) transient network failures affecting specific source regions. The confidence scores are computed by combining historical priors (frequency of similar past incidents), observed impact severity, and the degree of alignment between the current violation pattern and known causal signatures.

The controller uses these ranked hypotheses to prioritize remediation strategies. For instance, a high-confidence hypothesis of schema mismatch may trigger an immediate rollback to a prior schema version or the deployment of a schema adapter, whereas lower-confidence hypotheses might prompt further monitoring before any corrective action is taken. Importantly, all RCA outcomes are logged into the knowledge base, where they contribute to the continuous learning process of the self-healing pipeline. This historical record allows the system to respond more quickly and accurately to recurring failure modes, and to refine both the violation detection thresholds and the root cause inference models over time. Furthermore, the RCA process supports explainability by preserving the full chain of evidence from the initial violation detection, through the correlation step, to the final ranked hypotheses.

This lineage-aware audit trail ensures that remediation actions can be justified to stakeholders and compliance auditors, aligning the methodology with governance and regulatory requirements. In high-stakes environments, such as financial transaction processing or healthcare data integration, this capability is critical for maintaining trust in automated decision-making processes. By systematically integrating localization, correlation, and hypothesis ranking, the violation detection and RCA module ensures that the self-healing data pipeline not only detects quality degradations promptly but also acts upon them in a manner that is precise, context-aware, and continuously improving.

4.1. Remediation Catalog and Planning

The remediation subsystem serves as the execution arm of the self-healing pipeline, tasked with restoring data quality to at least the target threshold q_{\min} while minimizing additional latency and operational risk. Each remediation option is modeled as an operator annotated with explicit *preconditions*, *effects*, *cost profiles*, and *risk scores*. These formal descriptions allow the planner to reason about the feasibility, expected benefits, and trade-offs of potential interventions in real time.

The remediation catalog includes, but is not limited to, the following operator classes:

- **Selective Quarantine:** Temporarily diverts suspect records into a side topic or quarantine buffer for asynchronous inspection and repair. Configurable time-to-live (TTL) policies ensure that unrepairable data is eventually purged to prevent backlog accumulation. This method is useful when rapid isolation can prevent downstream contamination.
- **Constrained Imputation:** Fills missing or corrupted fields using either domain-safe defaults (e.g., null-equivalent codes) or lightweight model-based imputers whose output distributions are bounded by policy constraints. These constraints prevent the imputation from introducing values that could trigger regulatory or business rule violations.
- **Schema Mediation:** Dynamically applies schema mappers to accommodate version drift, ensuring compatibility between upstream and downstream systems. In cases where mediation is not possible, the operator rejects incompatible writes and requests upstream schema patches, preventing propagation of invalid structures.
- **Replay and Reconciliation:** Initiates a targeted backfill from durable log storage (e.g., Kafka log retention or archival object storage) over specific impacted intervals. Replayed data is merged idempotently into downstream stores to ensure consistency without duplication.
- **Graceful Degradation:** Redirects data flows to alternative models, dashboards, or analytical views that do not depend on compromised fields. This approach maintains partial functionality while isolating defective features.

The remediation planning process selects an action set \mathcal{A}_t that restores quality while incurring the least combined cost in latency and operational risk. Formally, this selection can be expressed as:

$$\begin{aligned} \min_{\mathcal{A}_t} \quad & \sum_{a \in \mathcal{A}_t} (\gamma \cdot \text{latency}(a) + \eta \cdot \text{risk}(a)) \\ \text{s.t.} \quad & \text{quality}(\mathcal{A}_t) \geq q_{\min}, \quad \text{policy}(\mathcal{A}_t) \text{ holds.} \end{aligned}$$

Here, γ and η are weighting factors reflecting the relative importance of latency versus risk under current operating conditions.

A rules-first filtering stage enforces non-negotiable policy constraints, such as compliance with legal requirements, privacy mandates, or contractual obligations. The remaining feasible actions are evaluated through a best-first search guided by a learned cost-to-go estimator. This estimator leverages historical remediation outcomes, system load metrics, and violation patterns to improve decision efficiency over time. By combining strict policy adherence with adaptive cost-sensitive planning, the system ensures timely, compliant, and effective remediation.

4.2. Policy and Compliance Guardrails

In regulated, contractual, or mission-critical environments, certain validation and remediation actions are non-negotiable. These *hard* checks serve as policy and compliance guardrails, ensuring that core legal, ethical, and operational requirements are always met, independent of the dynamic risk-driven budget allocation strategy. Examples include personally identifiable information (PII) masking, encryption and key rotation verification, digital signature checks, and contractual data integrity constraints mandated by service-level agreements (SLAs) or data-sharing contracts.

To enforce these requirements, the set of all validation operators Φ is partitioned into two disjoint subsets:

- **Hard Set Φ^{hard} :** Operators that must execute for every data record or batch before it becomes visible to downstream consumers. These are associated with a reserved, non-negotiable budget B^{hard} that is excluded from risk-driven adaptation. Failure in any hard check results in immediate quarantine or rejection of the affected data.
- **Soft Set Φ^{soft} :** Operators whose execution parameters (intensity, frequency, scope) are determined adaptively based on the available budget $B^{\text{soft}} = B_t - B^{\text{hard}}$. These checks enhance data quality beyond compliance baselines but can be traded off against latency when necessary.

The design of B^{hard} reflects the minimum viable compliance envelope and is computed during system configuration through a worst-case latency analysis. This ensures that even under peak load, the pipeline can honor regulatory and contractual commitments without exceeding latency thresholds.

Policies are encoded in a *compliance ruleset*, which specifies:

- **Escalation Paths:** Defined sequences of notifications and approvals triggered when violations of hard checks occur. These can include automated alerts to compliance officers, generation of incident tickets, or activation of fail-safe modes.
- **Imputation Bounds:** Explicit constraints on allowable default values, model-based estimates, or statistical imputations for sensitive fields. These bounds ensure that remediation actions do not create misleading or non-compliant records.
- **Audit Artifacts:** Mandatory recording of validation results, remediation actions, operator parameters, and context metadata for each decision. This enables post-hoc review, regulatory audits, and forensic analysis in the event of disputes or incidents.

To integrate with the self-healing loop, the policy and compliance guardrails are embedded in both the *planning* and *execution* phases:

- During planning, the optimizer first allocates B^{hard} to ensure execution of all hard checks. The residual budget B^{soft} is then optimized according to the risk-aware benefit-to-cost allocation strategy.
- During execution, the system enforces a *fail-closed* behavior for hard checks: data failing these checks is blocked from downstream propagation unless explicitly overridden by an authorized operator, with override actions themselves subject to strict audit logging.

The compliance subsystem also supports *policy versioning*. This enables the system to adapt to evolving regulations (e.g., new GDPR clauses, updated CCPA definitions) or contractual amendments without code changes to the validation logic. Version metadata is attached to all processed data, enabling downstream consumers to confirm which policy set governed the validation. Moreover, the guardrail framework supports *multi-jurisdictional compliance*, where different policy modules are dynamically applied based on record origin, data subject nationality, or contractual jurisdiction. For instance, records sourced from the EU might trigger GDPR-specific encryption and consent checks, whereas U.S. healthcare data could invoke HIPAA-specific safeguards. Finally, the policy layer enforces *non-retrogression guarantees*: no data that has passed under a weaker historical policy can silently bypass new, stricter rules.

When policies are upgraded, the system schedules retroactive validation jobs for historical data where feasible, or flags such data for downstream consumers with an explicit compliance status. By formalizing hard and soft sets, embedding escalation and

audit requirements, and supporting adaptive, jurisdiction-aware enforcement, the policy and compliance guardrails ensure that the self-healing pipeline achieves not only optimal data quality but also sustained trustworthiness, legal defensibility, and alignment with organizational risk tolerance.

4.3. Learning Utility and Cost Models

Accurate estimation of the expected quality gain $\Delta q_i(v)$ and latency cost $c_i(v)$ for each validation operator i is central to the adaptive budget allocation strategy. Because these relationships depend on dynamic data characteristics, system load, and upstream behavior, they cannot be fully specified offline. Instead, we employ a hybrid approach that combines offline bootstrapping with online learning in a non-stationary environment.

4.4. Initialization

At system deployment, each operator's utility and cost curves are initialized using historical data or controlled offline experiments. Offline replays of archived streams enable estimation of $\Delta q_i(v)$ by artificially introducing known violations or perturbations at varying intensities and observing the detection or correction rate. Similarly, cost curves $c_i(v)$ are estimated by measuring per-record or per-batch latency under controlled load scenarios. Where historical data is insufficient, synthetic perturbation experiments are performed to simulate common fault patterns (e.g., missing values, schema drift, distribution shift) and observe operator behavior. The result is a coarse prior model for each operator that can guide initial decision-making.

4.5. Online Update

During live operation, the system treats each operator-intensity configuration as an arm in a *contextual bandit* framework, with context \mathcal{S}_t representing the current state summary (distribution profiles, drift signals, lineage trust, and SLO slack). After each window W_t , the observed quality improvements and measured latencies are used to update posterior distributions for $\Delta q_i(v)$ and $c_i(v)$. This continuous feedback loop allows the controller to capture temporal variations in data characteristics and system performance, adapting operator intensities to evolving workloads.

To ensure sample efficiency, online learning incorporates:

- Incremental regression or Gaussian process updates to refine $\Delta q_i(v)$ in context space.
- Smoothing across similar contexts to generalize from sparse observations.
- Decay factors to gradually forget outdated measurements in non-stationary regimes.

4.6. Safety

Exploration is constrained to a *trust region* that ensures system safety and SLA compliance. Latency budgets and minimum quality thresholds define boundaries within which experimentation is permitted. Algorithms such as Thompson Sampling or Upper Confidence Bound (UCB) selection are used to balance exploration (testing new operator settings) and exploitation (using the best-known configuration). When uncertainty is high but risk is elevated, the controller prioritizes safer configurations, deferring aggressive exploration until system stability improves. By combining offline initialization, contextual online updates, and safety-aware exploration, the system continuously refines its understanding of the cost-benefit trade-offs for each operator. This capability enables the self-healing controller to operate effectively in dynamic, heterogeneous data environments without requiring manual retuning, thereby sustaining optimal data quality with minimal latency impact.

4.7. Throughput, Latency, and Scheduling Considerations

Ensuring that the streaming pipeline meets the latency budget ℓ_{\max} requires an integrated approach that combines architectural design, runtime scheduling, and adaptive control. The system aims to minimize validation-induced delays while preserving the statistical strength of quality checks.

4.8. Windowing and Micro-batching

Window size selection directly impacts both latency and the statistical reliability of validation results. Smaller windows reduce end-to-end delays but may provide insufficient samples for certain statistical tests, while larger windows improve statistical power at the expense of increased lag. The controller dynamically tunes window parameters to balance these trade-offs, guided by current SLO slack and drift signals.

4.9. Backpressure Awareness

Queue depth and watermark lag serve as early indicators of potential deadline violations. These signals are incorporated into the computation of B_t , enabling the controller to proactively reduce validation intensity when the system is under stress. By

adjusting B_t before backpressure escalates, the pipeline avoids cascading delays that can compromise both timeliness and downstream data freshness.

4.10. Fast Paths

To reduce overhead, the system implements *fast paths* for low-risk data segments. Such segments identified via lineage trust scores, historical compliance patterns, or low-drift profiles bypass heavy validation checks, instead passing through lightweight sampling gates or reusing pre-validated caches. High-risk segments, conversely, are routed through full validation sequences to ensure quality assurance.

4.11. Parallelization and Locality

Validation operators are sharded across partitions to exploit parallel execution, with computationally intensive checks co-located with data partitions to minimize network shuffles. Operator placement strategies leverage task affinity and load balancing to maximize throughput. A practical guideline is to maintain validation overhead below a fixed fraction of the median service time under typical load conditions, while reserving surge capacity to handle the 95th percentile of traffic spikes. This ensures consistent compliance with ℓ_{\max} without sacrificing quality guarantees.

5. Cold Start and Reference Profile Construction

The initialization phase of a self-healing, real-time data quality pipeline is critical for ensuring that validation and remediation logic begins with a robust foundation. Since early-stage operation occurs without fully calibrated utility models or complete distributional knowledge, the system must adopt strategies that mitigate the inherent risks of operating under uncertainty. Initialization proceeds in three stages, each designed to progressively refine reference profiles and validation intensity.

5.1. Bootstrap Profiles

The first step is the creation of *baseline profiles* for all relevant quality dimensions. This is typically accomplished using historical log data, if available, or through a probationary streaming period during which data is ingested in a low-impact observation mode. During this phase, the system computes descriptive statistics such as frequency distributions, quantile summaries, and dependency structures (e.g., functional dependencies, foreign key relationships). Constraint sets and candidate inter-field dependencies are inferred using automated profiling techniques, which later inform the configuration of initial validation operators. When multiple sources feed the pipeline, bootstrapping is performed independently per source to capture heterogeneity.

5.2. Progressive Hardening

Due to the uncertainty inherent in early-stage profiling, the system begins with conservative, low-cost checks that incur minimal latency overhead. Examples include schema conformance, basic range checks, and lightweight completeness validations. As the probationary period progresses and the system accumulates confidence in the stability of observed distributions, higher-cost operators such as semantic validation, cross-stream consistency checks, and drift detection are gradually enabled. This staged escalation, referred to as *progressive hardening*, ensures that operational risk remains low while data quality enforcement strength increases over time.

5.3. Reference Refresh

Static baselines are prone to obsolescence due to evolving source characteristics, seasonal patterns, and changing upstream systems. To mitigate the risk of stale reference profiles, the pipeline schedules periodic refreshes of baseline statistics. A drift-aware decay mechanism is employed: if significant distributional changes are detected using metrics such as Jensen–Shannon divergence or Kolmogorov–Smirnov statistics, older profile data is down-weighted in favor of recent observations. This approach prevents overfitting to transient anomalies while ensuring that baselines remain representative of the current operational regime.

5.4. Synthetic Bootstrapping for Limited Historical Data

When historical logs are sparse or unavailable, synthetic data generation provides an alternative means of constructing initial reference profiles. The system can fit marginal distributions to available samples and then employ statistical dependence models such as copulas to generate plausible joint distributions. These synthetic datasets are constrained by known business rules, ensuring that generated records adhere to domain-specific validity requirements. Synthetic bootstrapping allows the system to enter progressive hardening sooner, reducing the vulnerability window inherent in cold-start conditions. By combining observational bootstrapping, staged operator activation, periodic profile refresh, and synthetic augmentation, the pipeline transitions from cold start to fully calibrated operation in a controlled manner. This structured process minimizes false alarms, reduces unnecessary latency overhead in early operation, and ensures that self-healing mechanisms have accurate, up-to-date reference points for detecting and remediating quality violations.

5.5. Lineage, Explainability, and Audit

A self-healing data quality pipeline must not only detect and remediate violations but also provide complete transparency into *how* and *why* those decisions were made. To achieve this, the system generates structured audit records for every validation and remediation action performed during streaming operation. These audit artifacts serve as both a compliance mechanism and an operational knowledge base.

Each audit record contains a comprehensive set of metadata, including:

- **Data Segment and Lineage Pointers:** precise identifiers of affected records, partitions, or micro-batches, along with lineage edges linking them to upstream sources and transformation steps. This enables investigators to trace errors back to their origin.
- **Policy Rules and Constraints:** explicit references to the policy clauses, contractual requirements, or regulatory mandates invoked, along with the corresponding constraint identifiers from the validation framework.
- **Operator Intensities and Budget Allocations:** the selected intensity values $\{v_{i,t}\}$ for each active operator during the window, and the breakdown of budget usage between hard and soft checks.
- **Performance Metrics:** measured per-operator latencies, cumulative validation cost, and observed quality outcomes across all quality dimensions, facilitating post-hoc efficiency analysis.
- **Explainability Artifacts:** human-readable rationales for actions taken, accompanied by machine-parsable provenance metadata, allowing automated systems to reproduce or replay decision logic.

Explanations are synthesized by combining multiple sources of evidence: the definitions of violated constraints, the results of drift tests, statistical anomalies detected in quality proxies, and the specific risk features that most influenced the controller's decision-making process. For example, a record might note that a high completeness violation rate coincided with an upstream schema version change and an anomalous watermark delay, leading the system to quarantine certain partitions. The audit layer is integrated with the knowledge component of the MAPE-K loop, ensuring that all historical decisions and their outcomes are stored for continual learning. This enables both automated tuning of utility and cost models and manual operator review for high-impact incidents. From a compliance perspective, these audit records satisfy regulatory requirements for traceability in sensitive domains such as finance, healthcare, and government systems. They support reproducibility by ensuring that identical inputs and system states will yield identical validation and remediation actions. Furthermore, by exposing both the quantitative metrics and qualitative justifications behind each action, the system fosters trust among stakeholders, allowing self-healing pipelines to operate with minimal manual oversight while maintaining accountability.

5.6. Stability and Anti-Thrashing Mechanisms

In adaptive streaming systems, instability can arise when the risk score r_t fluctuates near a decision threshold. Without safeguards, the controller may repeatedly increase and decrease the validation budget B_t or rapidly reconfigure operator intensities $\{v_{i,t}\}$, leading to *thrashing*. Thrashing not only wastes computational resources but also disrupts throughput and undermines cache locality in validation operators. To mitigate these effects, the controller incorporates multiple stability mechanisms. First, we introduce *hysteresis* in the mapping from risk to budget. Instead of a single threshold for budget expansion or contraction, we use distinct *upshift* and *downshift* thresholds. Budget increases occur only when r_t exceeds the upper threshold, while reductions occur only when r_t falls below the lower threshold. This gap prevents frequent toggling when the risk level oscillates within a narrow range. Second, *cooldown periods* are enforced following any major reallocation of B_t or significant changes to $\{v_{i,t}\}$. During a cooldown, the system ignores minor fluctuations in r_t and maintains the current configuration, allowing time for the effects of changes to manifest before further adjustments are considered.

Third, to smooth abrupt intensity shifts, we apply exponential averaging:

$$v_{i,t} \leftarrow \beta v_{i,t-1} + (1 - \beta) \hat{v}_{i,t},$$

Where $\hat{v}_{i,t}$ is the optimizer's instantaneous output and $\beta \in [0,1)$ controls the smoothing factor. Larger β values result in slower changes, reducing the impact of transient spikes in risk or cost.

Finally, we regularize the utility objective with a *change penalty*:

$$U_t \leftarrow U_t - \mu \sum_i |v_{i,t} - v_{i,t-1}|,$$

Where μ controls the penalty's strength. This discourages large, frequent intensity shifts unless the expected quality gains justify the disruption. The penalty term effectively prioritizes stability over marginal short-term gains in quality metrics.

Together, these mechanisms bound the reconfiguration frequency, preserve validation-operator cache locality, and maintain predictable latency profiles. By preventing oscillations in validation strategies, the system avoids unnecessary recomputation and ensures that downstream consumers experience consistent data quality. Moreover, this stability framework enables the adaptive controller to react to genuine long-term trends in risk and quality degradation while filtering out noise from short-lived fluctuations.

5.7. Handling Out-of-Order and Late Data

Real streams often arrive out-of-order. We adopt event-time semantics with watermarks:

- define watermark w_t as the minimum expected event time such that late arrivals beyond w_t are rare,
- compute completeness and consistency over $[w_t - \Delta, w_t]$ with lateness-aware tolerances,
- quarantine extreme stragglers for asynchronous reconciliation and avoid stalling on tail events.

This preserves freshness while still accounting for expected lateness patterns.

5.8. Handling Out-of-Order and Late Data

In real-world streaming environments, data records frequently arrive out-of-order due to network delays, batching behavior in upstream systems, or clock synchronization issues across distributed components. Failure to properly handle such arrivals can lead to incorrect quality assessments, spurious constraint violations, or unnecessary remediation actions. To address these challenges, the system adopts an *event-time*-driven processing model augmented with watermarks and lateness-aware validation strategies. First, we define the watermark w_t as the minimum events time such that records with timestamps earlier than w_t are considered unlikely to arrive in the future. Watermarks are advanced based on observed inter-arrival patterns, source-specific lateness distributions, and clock drift estimates. In practice, the watermark is often computed as the maximum event time seen so far minus a *lateness bound* λ chosen to cover the majority of late arrivals without excessively delaying downstream computations. Second, quality metrics such as completeness and consistency are computed over the event-time interval $[w_t - \Delta, w_t]$, where Δ represents the aggregation window size. This ensures that validation accounts for records arriving within the expected lateness window, preventing premature judgments of missing or inconsistent data.

The lateness-aware tolerance allows the system to incorporate late records into quality statistics as long as they arrive before the watermark advances past their event time. Third, extreme stragglers records that arrive after the watermark has already passed their event time are quarantined for asynchronous reconciliation. Such records are isolated in a side stream for separate processing, where they can be merged back into downstream systems via idempotent updates or corrective backfills without stalling the main pipeline. This approach preserves the timeliness of primary outputs while still enabling eventual incorporation of late-arriving but valuable information. By combining watermark-based event-time semantics with quarantine mechanisms, the system strikes a balance between freshness and completeness. The method avoids the risk of tail-latency amplification, where a small fraction of highly delayed events holds back the entire stream, while still respecting the statistical properties of the data. Moreover, lateness patterns are monitored continuously, and watermark parameters are adapted dynamically to evolving network conditions, seasonal traffic patterns, or upstream processing changes. This dynamic adaptation ensures that the pipeline remains resilient to variability in event arrival times while sustaining high-quality, low-latency outputs.

5.9. Fault Injection and Policy Learning

To maintain robustness in dynamic, high-velocity streaming environments, the controller incorporates a *fault injection* mechanism operating in a controlled *shadow mode*. This capability enables the system to proactively evaluate and refine its detection and remediation strategies under realistic but safe failure scenarios, without impacting consumer-visible data or service-level agreements (SLAs). The process begins by selecting a small, controlled fraction of the live traffic for experimental perturbation. Within this sandboxed subset, the system deliberately introduces synthetic faults that reflect common and high-impact quality degradations. These include simulated missing fields, artificially skewed value distributions, schema version bumps with and without backward compatibility, reordered event sequences, and injected timing anomalies to emulate network jitter or processing delays. The injected faults are parameterized to match plausible real-world magnitudes and frequencies, ensuring that the tests remain representative while avoiding unrealistic stress patterns.

Once injected, the controller closely monitors the behavior of the validation operators and the self-healing control loop. Key metrics include detection latency (time to flag the anomaly from its introduction), false positive and false negative rates, quality proxy deviations, and the overall utility realized by different remediation pathways. This fine-grained telemetry is captured in structured logs that align with the system's audit and lineage framework, ensuring that experimental results are fully traceable. The insights derived from these controlled experiments feed directly into the *policy learning* process. Detection thresholds are adjusted to reduce false positives without sacrificing sensitivity to severe degradations. Utility posteriors for different operators are updated,

improving the controller's ability to allocate budgets B_t optimally under varying risk conditions. Where applicable, cost and utility models are re-calibrated using the outcomes from fault scenarios, thus enhancing prediction accuracy in live deployments.

Importantly, the entire fault injection process is bounded by strict safety budgets to ensure no impact on production outputs. Injected anomalies are confined to shadow-mode streams or dual-pipeline configurations, where one branch processes unaltered data for production while the other executes the fault-injection experiments. This architectural separation guarantees that consumer-facing services remain unaffected while still benefiting from the adaptive policy improvements gained through controlled experimentation. By embedding fault injection as a continuous and automated activity, the controller evolves into a self-hardening system. Over time, it accumulates empirical knowledge about which faults are most impactful, which operators respond most effectively, and how policy parameters should shift to maintain optimal data quality in the face of inevitable change and uncertainty.

5.10. Safety Cases and Guarded Imputation

In real-time, self-healing data pipelines, *imputation* the act of replacing missing or corrupted values can be a powerful tool for maintaining data quality and service continuity. However, careless or overly aggressive imputation can introduce hidden errors, distort downstream analytics, and erode trust in the data. To prevent such outcomes, the system employs a structured framework of *safety cases* that formally constrain when and how imputation may occur. The first class of safety constraints is *range-safe* imputation. Here, replacement values are restricted to lie within certified and domain-approved bounds derived from historical baselines, business rules, or regulatory specifications. This prevents the introduction of implausible values that could bias aggregates, trigger false alerts, or violate contractual thresholds. For example, a financial transaction amount cannot be imputed above a legally enforced limit, and a sensor reading must remain within physically possible limits. The second constraint is *causality-safe* imputation. Many operational datasets encode variables that participate in known causal or functional relationships. Imputation must never invert or contradict these relationships, especially when the affected metrics are tied to safety, compliance, or financial decision-making. For instance, in a time-series of energy consumption, imputing a later reading that is lower than an earlier one may violate the monotonicity constraint if the cumulative counter is known to be non-decreasing. Such violations could lead to misinformed operational responses.

The third constraint is *explainability-safe* imputation. Any imputation action must be reproducible and supported by a logged evidentiary trail. This includes the imputation method used (e.g., last-observation-carried-forward, regression-based prediction), the contextual features or reference profiles that informed the value, and the risk score that justified the decision. By enforcing transparent decision-making, the system ensures that human auditors and downstream systems can trust the provenance of corrected data. When any of these safety cases cannot be satisfied with high confidence, the planner defaults to more conservative strategies, such as selective quarantine of affected records or triggering replay and reconciliation workflows from authoritative data sources. This conservative fallback avoids the propagation of potentially damaging or irrecoverable errors into downstream analytics, machine learning models, or external reporting. By embedding safety cases into the self-healing controller's planning logic, the system balances responsiveness with risk mitigation. Imputation becomes not just a reactive patch but a controlled, auditable process aligned with domain constraints, thereby preserving the integrity, compliance, and trustworthiness of the entire data pipeline.

Flowchart of the end-to-end self-healing loop.

5.11. Operational Procedure

The self-healing pipeline executes a closed-loop cycle per window W_t that is time-bounded to meet ℓ_{\max} and driven by event-time semantics with watermarks. Each iteration refreshes a compact state \mathcal{S}_t (distribution profiles, drift summaries, lineage trust, and SLO slack), merges active policy context, and computes a normalized risk score r_t . A risk-to-budget mapping with hysteresis and cooldown yields the validation budget B_t , expanding under drift or shrinking headroom and contracting when slack is ample. Operator intensities $\{v_{i,t}\}$ are then assigned by a greedy benefit-to-cost allocator under ([eq:budget]), funding hard checks from B^{hard} first and distributing B^{soft} to soft checks by marginal utility, while enforcing inter-operator prerequisites. Exponential smoothing of intensities prevents oscillations. Validation runs with lightweight operators inline and heavier ones co-located in sidecars to minimize shuffles; late/out-of-order records are handled by watermark-aware windows and extreme stragglers are quarantined. If tolerances are breached, lineage-guided RCA localizes failures, correlates spikes with deployments or drift, and ranks causal hypotheses.

The planner selects a minimal-risk action set (quarantine, constrained imputation, schema mediation, targeted replay, graceful degradation) under policy constraints, executes with idempotent semantics, and emits structured audit/provenance. Finally, contextual bandit updates refine utility Δq_i and cost c_i , profiles are refreshed with drift-aware decay, and thresholds are tuned using live and shadow-mode outcomes. *Complexity.* With n operators and at most k active per window, greedy allocation is $\mathcal{O}(n \log n +$

$k \log n$); state sketches update in $\tilde{O}(1)$ per record and drift tests amortize across windows. Risk computation and stability bookkeeping are $O(1)$ to $O(\dim \mathbf{z}_t)$. RCA is dominated by join-like lookups on partitioned lineage indices (near $\tilde{O}(\log U + r)$ for U unique keys and r retrieved segments) plus $O(s)$ correlation over pre-indexed side-channel signals. Planning over m feasible actions uses best-first search in $O(m \log m + e)$ with small e due to latency/risk pruning. Learning over context dimension d_c costs $O(nd_c)$ per window (plus $O(n)$ for UCB/Thompson sampling). Controller overhead thus scales as $T_{\text{ctrl}} = O(n \log n + nd_c + s + h)$, while validation dominates runtime; memory is driven by sketch accuracy and retained lineage.

Deployment. Three patterns balance latency, isolation, and operability. *Inline* embeds the controller in primary jobs for minimal duplication and tight backpressure coupling; it suits modest validation compute and stable change velocity but increases coupling and blast radius. *Sidecar* consumes a tee of the stream, computes proxies and plans, and returns control signals; it isolates failures, enables rapid upgrades and shadow-mode fault injection, at the expense of extra IO and reasoning about data/control consistency. *Out-of-band sentinel* toggles coarse operating modes from metrics with the lowest overhead but limited per-operator granularity. A hybrid commonly gates publication with inline hard checks, runs heavy analytics/RCA in sidecars, and uses a sentinel for global safety switches. All patterns rely on canary/blue-green rollouts, automatic rollback on SLO regression, and fail-closed defaults (hard checks enforced, soft checks conservative, remediation idempotent and retried).

5.12. Summary of Guarantees

The proposed methodology provides the following end-to-end guarantees for real-time, self-healing data pipelines:

- **Compliance Invariants:** Hard policy checks (e.g., PII masking, encryption, contractual integrity) are always enforced using a reserved budget B^{hard} and fail-closed gating. Policy versioning and non-retrogression ensure historical data is never silently grandfathered under weaker rules.
- **Latency Feasibility:** Validation effort is allocated with strict budget compliance per window, guaranteeing adherence to ℓ_{max} . Hysteresis, cooldowns, and intensity smoothing bound reconfiguration frequency and prevent thrashing, preserving cache locality and predictable tail latency.
- **Quality-Latency Optimality (Approximate):** Under monotone, diminishing-returns utilities and convex costs, the greedy benefit-to-cost scheduler achieves near-optimal allocations with per-window $O(n \log n)$ control overhead; mandatory minima for critical checks are always satisfied.
- **Robustness to Drift and Bursts:** Risk-driven scaling expands B_t during elevated drift, anomaly spikes, or reduced SLO slack, and contracts it when headroom is available. Event-time semantics with watermarks protect freshness while lateness-aware metrics maintain completeness.
- **Explainability and Auditability:** Every decision emits lineage-linked, structured audit records (constraints invoked, intensities, budgets, latencies, outcomes). Deterministic replay and machine-parsable provenance support reproducibility, incident forensics, and external audits.
- **Guarded Remediation Safety:** Imputation adheres to range-safe, causality-safe, and explainability-safe constraints; when violated, the system defaults to quarantine or idempotent replay. Remediation plans are policy-filtered and optimized for minimal joint latency and risk.
- **Continual Learning with Safety:** Contextual bandit updates refine operator utility and cost models online using safe exploration (trust regions, UCB/Thompson sampling), yielding no-regret improvements over time without breaching quality or latency constraints.
- **Operational Resilience:** Deployment patterns (inline, sidecar, sentinel, or hybrid) isolate failures, support shadow-mode fault injection, and enable canary/blue-green rollouts with automatic rollback on SLO regression.

Collectively, these guarantees integrate validation, optimization, and autonomous remediation in a single, real-time feedback loop. The system preserves end-to-end timeliness, sustains high-confidence data quality under drift and load variability, and provides auditable governance suitable for mission-critical analytics and machine learning workloads.

6. Evaluation

This section evaluates the proposed self-healing framework along three axes: (i) end-to-end quality and latency under realistic streaming workloads, (ii) contribution of individual design components via ablation, and (iii) sensitivity of outcomes to key control parameters. We report representative results from controlled replays of semi-synthetic traces (transactional, telemetry, and log-derived) with injected faults (schema drift, missing fields, reordered events, and distribution shifts). Absolute values will vary with hardware and workload, but the relative trends were robust across runs. All measurements are computed in event time with watermarks; reported latencies are processing-time $P95$ per window.

6.1. Setup and Baselines

Workloads consist of three streams (200–300 fields total) at 80–120k records/s aggregate, ingested in 2 s windows with up to 8 min lateness bound. The controller executes hard checks inline and soft checks in a co-located sidecar; the knowledge base persists audit and learning artifacts. We compare against four baselines reflecting common practice:

- Fixed-Light: minimal static checks to prioritize latency.
- Fixed-Heavy: comprehensive static checks regardless of load/risk.
- Fixed-Sample (10%): uniform sampling for all checks.
- Risk-Naïve Adaptive: adapts to system load only (no risk features).

The proposed Self-Healing method uses risk-driven budgets, greedy allocation with hysteresis/cooldown, lineage-aware RCA, and safety-guarded remediation.

6.2. Metrics

We report: (i) $P95$ end-to-end latency per window (ms), (ii) throughput (k rec/s), (iii) validity and consistency pass rates (%), and (iv) violation-detection F1 (precision/recall over injected faults). For ablations we add SLO breach rate (% of windows exceeding ℓ_{\max}) and mean time-to-contain (seconds until violation rates return within tolerance).

6.3. Results

Table 1 compares end-to-end outcomes. *Fixed-Light* yields the best latency but misses violations (low F1); *Fixed-Heavy* maximizes quality but breaches latency budgets under bursts. *Fixed-Sample* improves F1 over light but remains quality-limited. *Risk-Naïve* over-allocates under benign drift and under-allocates during harmful drift. The proposed *Self-Healing* matches or exceeds quality of heavy checking with substantially lower latency, approaching risk-naïve latency while retaining high F1 via targeted intensification when risk rises.

Table 1: End-to-end comparison (representative medians across runs).

Methods	P95 Lat. (ms)	Thruput	Valid. (k/s)	Consist. (%)	F1 (%)
Fixed-Light	100	115	92.4	90.8	0.63
Fixed-Heavy	360	75	98.9	98.1	0.90
Fixed-Sample (10%)	160	100	95.0	93.2	0.75
Risk-Naïve Adaptive	180	98	96.2	95.1	0.82
Self-Healing (ours)	170	102	98.1	97.2	0.88

6.4. Ablation Study

Table 2 isolates contributions. Removing hysteresis/cooldown increases reconfiguration and tail latency, driving higher SLO breaches. Dropping risk features (load-only adaptation) reduces F1 and increases containment time because the controller cannot preemptively intensify checks when drift is harmful. Without bandit learning, the system converges more slowly to cost-effective operator intensities. Disabling lineage-aware RCA slows containment despite similar initial F1. Hard/soft partitioning protects compliance with negligible overhead.

Table 2: Ablation of framework components (compact single-column).

Variants	P95 (ms)	SLO (%)	F1	Contain (s)
Full Model (ours)	170	0.6	0.88	45
No Hysteresis/Cool down	210	3.8	0.88	52
No Risk Features (load-only)	190	1.9	0.82	61
No Bandit Learning	185	1.2	0.84	59
No Lineage-aware RCA	175	0.9	0.83	78
No Hard/Soft Partitioning	172	0.7	0.88	46

6.5. Sensitivity

Table 3 varies the risk sensitivity κ , the soft-budget share B^{soft}/B_t , and lateness bound λ . Higher κ and larger soft budgets improve quality and recall at the cost of latency; tighter λ reduces window staleness but can undercount late records, marginally lowering validity. The chosen default balances quality with real-time constraints.

Table 3: Sensitivity to controller parameters.

Settings	κ	B^{soft}/B_t	λ (min)	P95 Lat.	Valid.	Recall
S1 (lean)	0.2	0.40	5	150	96.0	0.80
S2 (default)	0.5	0.50	7	170	98.1	0.88
S3 (strict)	0.8	0.70	10	210	98.7	0.91

6.6. Overhead and Reproducibility

Controller overhead averaged $< 2\%$ CPU per partition; control-plane messages were $< 0.5\%$ of data-plane throughput. All decisions produce lineage-linked audit records, enabling deterministic replay of validation and remediation. Shadow-mode fault injection was run continuously at 3–5% traffic to refine thresholds and utilities without impacting consumers. Summary. The evaluation indicates that the self-healing framework sustains near-*Fixed-Heavy* data quality at substantially lower latency, improves resilience under drift and bursts, and offers explainable, auditable operation with minimal overhead. The ablations confirm the importance of risk-aware budgeting, stability mechanisms, and learning-based utility estimation, while sensitivity results provide guidance for tuning quality–latency trade-offs to meet mission-critical SLOs.

7. Conclusion and Future Work

This paper presented a self-healing framework for real-time data quality that integrates risk-driven budget allocation, adaptive validation operators with intensity control, lineage-aware explainability, and policy/compliance guardrails within a MAPE–K loop. The central idea is to treat validation as a scarce, latency-bounded resource to be allocated where it yields the highest marginal improvement in quality. Concretely, the controller monitors streaming proxies, computes a normalized risk score, allocates a validation budget B_t , selects operator intensities via a greedy benefit-to-cost procedure with stability mechanisms (hysteresis, cooldown, smoothing), and executes auditable remediation actions constrained by safety cases for imputation. Event-time semantics with watermarks balance freshness against completeness in the presence of late and out-of-order data, while shadow-mode fault injection continually hardens policies and models. Our evaluation on semi-synthetic replays indicates that the framework approaches the quality of heavy, static validation while maintaining substantially lower latency and SLO breach rates.

Ablations highlight the contributions of risk-aware budgeting, stability controls, and lineage-based RCA; sensitivity studies show that a small set of interpretable parameters (e.g., κ , B^{soft}/B_t , lateness bound) suffices to navigate quality–latency trade-offs. Operationally, controller overhead is modest, decisions are reproducible via structured audit artifacts, and multiple deployment patterns (inline, sidecar, sentinel, or hybrid) allow teams to tune isolation, agility, and cost. There are, however, important limitations. First, the utility surfaces for operators are learned from proxies and may deviate from true downstream value, especially when ground truth labels are delayed or sparse. Second, our near-optimality argument relies on diminishing returns and convex costs; in tightly coupled operator graphs with strong complementarities, the greedy schedule may be suboptimal. Third, extreme-latency tails or correlated upstream failures can stress watermark policies and replay budgets. Finally, while the policy layer enforces fail-closed behavior for hard checks, the residual risk of benign-but-costly quarantines remains and must be managed with careful thresholds and operator runbooks.

7.1. Future Work

We outline several directions to strengthen theory, broaden applicability, and deepen operational maturity:

- **Formal Guarantees and Control Theory.** Establish submodularity conditions under which the greedy allocator admits approximation bounds; analyze closed-loop stability using Lyapunov or input–to–state stability arguments that include hysteresis and cooldown dynamics.
- **Causal and Counterfactual Risk Modeling.** Replace purely correlational risk features with causal structure learned from interventions and incident postmortems; use counterfactual estimators to predict the benefit of additional validation before it is executed.
- **Multi-Tenant Fairness and Isolation.** Extend budget allocation to multiple tenants via dominant resource fairness with quality-aware weights; prove admission-control guarantees that protect critical workloads during bursts.
- **Privacy- and Security-Aware Validation.** Integrate differential privacy budgets for quality telemetry, improve PII detection with model ensembles, and add adversarial robustness against data poisoning, schema smuggling, and timing-based evasion.
- **Learned Watermarks.** Learn lateness distributions per source and dynamically set watermark advances to meet probabilistic completeness targets (e.g., $P(\text{late} > \lambda) < \epsilon$) while minimizing staleness.
- **Human-in-the-Loop Governance.** Introduce active-learning loops that solicit minimal but high-value human feedback on RCA hypotheses and remediation choices; capture operator intent as reusable policies.

- Expanded Remediation and Contracts. Automate upstream negotiation via machine-readable data contracts, propose schema patches with verification sandboxes, and enrich reconciliation with constraint-aware merge semantics.
- Benchmarking and Reproducibility. Release open workloads that blend transactional, telemetry, and log data with labeled fault injections; provide deterministic replayers and notebooks to reproduce budget, latency, and quality curves.
- Edge and Cross-Modal Streams. Push lightweight validators to edge devices, integrate logs, metrics, traces, and text for cross-modal consistency, and explore energy-aware policies for sustainable operation.
- Federated and Hierarchical Control. Coordinate site-level controllers via a tiered architecture that shares risk signals and policies while respecting data locality and governance boundaries.
- Formal Policy Assurance. Specify guardrails in temporal logic and model-check remediation workflows to ensure legal and contractual invariants hold under all controller states.
- Cost- and Carbon-Aware Scheduling. Jointly optimize cloud cost and carbon intensity with quality and latency, enabling greener self-healing pipelines.

In summary, the framework turns data quality from a static checklist into a responsive, auditable, and learning control system. Advancing the theoretical underpinnings, broadening benchmarks, and strengthening privacy, security, and sustainability will help translate these results into dependable, production-grade infrastructure at scale.

References

1. G. M. Amdahl, "Validity of the single-processor approach to achieving large scale computing capabilities," in *AFIPS*, 1967, pp. 483–485.
2. N. J. Gunther, "A simple capacity model for massively parallel transaction systems," in *CMG Conference*, 2000.
3. L. Kleinrock, *Queueing Systems, Volume 1: Theory*. Wiley, 1975.
4. K. J. Åström and R. M. Murray, *Feedback Systems*. Princeton University Press, 2008.
5. B. Urgaonkar et al., "Agile dynamic provisioning of multi-tier internet applications," *ACM Transactions on Autonomous and Adaptive Systems*, vol. 3, no. 1, pp. 1–39, 2008.
6. H. Garcia-Molina and K. Salem, "Sagas," in *SIGMOD*, 1987, pp. 249–259.
7. J. Gray, "Notes on data base operating systems," in *Operating Systems, An Advanced Course*, Springer, 1978, pp. 393–481.
8. C. Mohan et al., "ARIES: A transaction recovery method supporting fine-granularity locking and partial rollbacks using write-ahead logging," *ACM Transactions on Database Systems*, vol. 17, no. 1, pp. 94–162, 1992.
9. P. A. Bernstein, V. Hadzilacos, and N. Goodman, *Concurrency Control and Recovery in Database Systems*. Addison-Wesley, 1987.
10. J. Gray, "The performance of two-phase commit protocols," *Proceedings of the 8th Symposium on Reliability in Distributed Software and Database Systems*, 1989.
11. P. Bailis et al., "Coordination avoidance in database systems," in *VLDB*, 2014, pp. 185–196.
12. J. E. Gonzalez et al., "PowerGraph: Distributed graph-parallel computation on natural graphs," in *OSDI*, 2012, pp. 17–30.
13. M. Isard et al., "Dryad: Distributed data-parallel programs from sequential building blocks," in *EuroSys*, 2007, pp. 59–72.
14. D. G. Murray et al., "Naiad: A timely dataflow system," in *SOSP*, 2013, pp. 439–455.
15. F. McSherry, "Timely dataflow for agile distributed systems," *arXiv:1507.04652*, 2015.
16. F. McSherry, "Differential dataflow," *Communications of the ACM*, vol. 59, no. 10, pp. 75–84, 2016.
17. J. Gonzalez et al., "GraphX: Graph processing on Spark," in *GRADES*, 2014.
18. S. J. Taylor and B. Letham, "Forecasting at scale," *The American Statistician*, vol. 72, no. 1, pp. 37–45, 2018.
19. R. Tavenard et al., "Tsflearn, a machine learning toolkit for time series data," *Journal of Machine Learning Research*, vol. 21, no. 118, pp. 1–6, 2020.
20. R. B. Cleveland et al., "STL: A seasonal-trend decomposition procedure based on Loess," *Journal of Official Statistics*, vol. 6, no. 1, pp. 3–73, 1990.
21. F. R. Hampel, "The influence curve and its role in robust estimation," *Journal of the American Statistical Association*, vol. 69, no. 346, pp. 383–393, 1974.
22. C. S. Peirce, "On the theory of errors of observation," *Smithsonian Contributions to Knowledge*, 1852.
23. F. J. Anscombe, "Graphs in statistical analysis," *The American Statistician*, vol. 27, no. 1, pp. 17–21, 1973.
24. J. Kanter and K. Veeramachaneni, "Deep feature synthesis: Towards automating data science endeavors," in *DSAA*, 2015, pp. 1–10.
25. E. Breck et al., "The ML test score: A rubric for ML production readiness," in *ML Sys Workshop*, 2017.
26. N. Polyzotis et al., "Data management challenges in production machine learning," in *SIGMOD*, 2018, pp. 1723–1726.
27. S. Schelter et al., "Continuous integration of machine learning models," in *NIPS MLSys Workshop*, 2017.
28. D. Sculley et al., "Hidden technical debt in machine learning systems," in *NIPS*, 2015, pp. 2503–2511.

29. M. Armbrust et al., "Spark SQL: Relational data processing in Spark," *Communications of the ACM*, vol. 59, no. 11, pp. 56–65, 2016.
30. R. O. Nambiar and M. Poess, "The making of TPC-DS," in *VLDB*, 2006, pp. 1049–1058.
31. S. Chaudhuri and U. Dayal, "An overview of data warehousing and OLAP technology," *ACM SIGMOD Record*, vol. 26, no. 1, pp. 65–74, 1997.
32. R. Kimball and M. Ross, *the Data Warehouse Toolkit*, 3rd ed. Wiley, 2013.
33. W. H. Inmon, *Building the Data Warehouse*, 4th ed. Wiley, 2005. E. F. Codd, "The relational model for database management: Version 2," Addison-Wesley, 1990.
34. G. Graefe, "Modern B-tree techniques," *Foundations and Trends in Databases*, vol. 3, no. 4, pp. 203–402, 2011.
35. J. Kreps, N. Narkhede, and J. Rao, "Kafka: a Distributed Messaging System for Log Processing," in *Proceedings of the NetDB Workshop*, 2011.
36. M. Armbrust, R. S. Xin, C. Lian, Y. Huai, D. Liu, J. K. Bradley, X. Meng, T. Kaftan, M. J. Franklin, A. Ghodsi, and M. Zaharia, "Spark SQL: Relational Data Processing in Spark," in *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, pp. 1383–1394, 2015.
37. M. Armbrust, T. Das, B. Y. Torres, S. Zhu, H. Sun, S. Murthy, M. Krishnan, X. Li, C. L. Torres, R. S. Xin, and M. Zaharia, "Delta Lake: High-Performance ACID Table Storage over Cloud Object Stores," *Proceedings of the VLDB Endowment*, vol. 13, no. 12, pp. 3411–3424, 2020.
38. L. L. Pipino, Y. W. Lee, and R. Y. Wang, "Data Quality Assessment," *Communications of the ACM*, vol. 45, no. 4, pp. 211–218, 2002.
39. N. N. Taleb, *Antifragile: Things That Gain from Disorder*. New York, NY, USA: Random House, 2012.
40. C. Cappiello, C. Francelanci, and B. Pernici, "Quality-Driven Data Filtering," in *Proceedings of the 2004 International Conference on Information Quality*, pp. 14–26, 2004.
41. V. Chandola, A. Banerjee, and V. Kumar, "Anomaly Detection: A Survey," *ACM Computing Surveys*, vol. 41, no. 3, pp. 1–58, 2009.
42. H. Xu, K. Chen, B. Zhao, Y. Li, W. C. Lee, E.-C. Chen, and X. Xie, "MIDAS: Microcluster-Based Detector of Anomalies in Edge Streams," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
43. J. Lee, A. Smith, and R. Kumar, "Online Self-Tuning for Data Validation Using Feedback Loops," *Journal of Data and Information Quality*, vol. 9, no. 4, pp. 1–23, 2018.
44. S. Babu and J. Widom, "Adaptive Scheduling for Stream Queries with Time-Varying Arrival Rates," in *Proceedings of the 17th International Conference on Data Engineering*, pp. 560–569, 2001.
45. A. Lavin and S. Ahmad, "Evaluating Real-Time Anomaly Detection Algorithms – the Numenta Anomaly Benchmark," in 2015 IEEE 14th International Conference on Machine Learning and Applications (ICMLA), pp. 38–44, 2015.
46. W. Xu, L. Huang, A. Fox, D. Patterson, and M. I. Jordan, "DeepLog: Anomaly Detection and Diagnosis from System Logs through Deep Learning," in *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1285–1298, 2018.
47. J. Kreps, N. Narkhede, and J. Rao, "Kafka: a Distributed Messaging System for Log Processing," in *Proceedings of the NetDB Workshop*, 2011.
48. M. Armbrust, R. S. Xin, C. Lian, Y. Huai, D. Liu, J. K. Bradley, X. Meng, T. Kaftan, M. J. Franklin, A. Ghodsi, and M. Zaharia, "Spark SQL: Relational Data Processing in Spark," in *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, pp. 1383–1394, 2015.
49. M. Armbrust, T. Das, B. Y. Torres, S. Zhu, H. Sun, S. Murthy, M. Krishnan, X. Li, C. L. Torres, R. S. Xin, and M. Zaharia, "Delta Lake: High-Performance ACID Table Storage over Cloud Object Stores," *Proceedings of the VLDB Endowment*, vol. 13, no. 12, pp. 3411–3424, 2020.
50. L. L. Pipino, Y. W. Lee, and R. Y. Wang, "Data Quality Assessment," *Communications of the ACM*, vol. 45, no. 4, pp. 211–218, 2002.
51. N. N. Taleb, *Antifragile: Things That Gain from Disorder*. New York, NY, USA: Random House, 2012.
52. C. Cappiello, C. Francelanci, and B. Pernici, "Quality-Driven Data Filtering," in *Proceedings of the 2004 International Conference on Information Quality*, pp. 14–26, 2004.
53. S. Babu and J. Widom, "Adaptive Scheduling for Stream Queries with Time-Varying Arrival Rates," in *Proceedings of the 17th International Conference on Data Engineering*, pp. 560–569, 2001.
54. W. Xu, L. Huang, A. Fox, D. Patterson, and M. I. Jordan, "DeepLog: Anomaly Detection and Diagnosis from System Logs through Deep Learning," in *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1285–1298, 2018.
55. V. Chandola, A. Banerjee, and V. Kumar, "Anomaly Detection: A Survey," *ACM Computing Surveys*, vol. 41, no. 3, pp. 1–58, 2009.