# Data Observability and Data Quality Automation: Building Self-Healing Data Pipelines

## Ramesh Betha

Independent Researcher
Holly Springs NC, US.
ramesh.betha@gmail.com

**Abstract:**
**Modern data architectures have become increasingly complex, creating new challenges in ensuring data quality and reliability. This paper explores the emerging field of data observability and quality automation frameworks that enable organizations to build self-healing data pipelines. We present a comprehensive analysis of current challenges in data quality management, examine the evolution of observability practices from DevOps to DataOps, and propose a reference architecture for implementing intelligent data quality systems. Through case studies and empirical evidence, we demonstrate how organizations can significantly reduce data downtime, accelerate issue resolution, and build greater trust in their data assets through automated detection, diagnosis, and remediation capabilities. The paper concludes with a roadmap for future developments in self-healing data systems and guidelines for implementation across various organizational contexts.**

**Keywords: data observability, data quality, self-healing systems, DataOps, automation, machine learning, data reliability engineering.**

## I. INTRODUCTION

Data has become the lifeblood of modern organizations, powering everything from operational decision-making to strategic planning and artificial intelligence initiatives. As data volumes grow exponentially and data architectures become increasingly complex, ensuring data quality and reliability has emerged as one of the most critical challenges facing data teams today [1]. Research by Gartner indicates that poor data quality costs organizations an average of $12.9 million annually, with impacts ranging from lost revenue to missed opportunities and diminished trust in data-driven decision-making [2].

Traditional approaches to data quality management—characterized by manual checks, reactive troubleshooting, and siloed responsibilities—have proven inadequate in modern data environments. As organizations adopt cloud-native architectures, microservices, and real-time data processing, the points of potential failure multiply while visibility into data pipelines decreases. This complexity has given rise to a new paradigm: data observability and automated quality management systems that can detect, diagnose, and even remediate issues with minimal human intervention.

This paper explores the emerging field of self-healing data pipelines, defined as data workflows that can autonomously monitor their own health, identify quality issues, determine root causes, and implement corrective actions. We examine the technological foundations, architectural patterns, and organizational practices that enable truly resilient data systems. Drawing on real-world implementations and research, we provide a blueprint for organizations seeking to move from reactive data firefighting to proactive data quality automation.

The remainder of this paper is organized as follows: Section II examines the current state of data quality challenges and the limitations of traditional approaches. Section III introduces the concept of data observability and its relationship to broader observability practices. Section IV presents a reference architecture for self-healing data pipelines. Section V explores implementation strategies and case studies.

Section VI discusses future directions and emerging technologies. Section VII provides concluding thoughts and practical recommendations.

## II. THE DATA QUALITY CRISIS: UNDERSTANDING THE CHALLENGE

### A. The Evolving Data Landscape

The data landscape has undergone a profound transformation over the past decade. Organizations have migrated from monolithic data warehouses to distributed architectures comprising data lakes, cloud-native services, and specialized processing engines. According to a survey by Dimensional Research, 92% of data engineers report that their data environments have become more complex in the past year [3]. This complexity manifests in several ways:

- *Volume and Velocity:* IDC projects that the global datasphere will grow from 64.2 zettabytes in 2020 to 180 zettabytes by 2025 [4]. This explosive growth is accompanied by increasing velocity, with more data being processed in real-time streams.
- *Variety and Heterogeneity:* Organizations routinely integrate data from dozens or hundreds of sources in different formats, structures, and levels of quality. Each integration point represents a potential failure point.
- *Distributed Ownership:* As data democratization efforts progress, responsibility for data creation, transformation, and consumption is distributed across the organization, complicating governance and accountability.
- *Technological Diversity:* Modern data stacks comprise numerous specialized tools and platforms—ETL tools, data catalogs, transformation frameworks, and analytics platforms—each with its own interfaces, semantics, and failure modes.

This evolving landscape has created what some practitioners call the "data quality gap"—the growing distance between the complexity of data systems and our ability to ensure their reliability.

### B. The Consequences of Poor Data Quality

The consequences of data quality issues extend far beyond technical inconvenience. Research by MIT Sloan indicates that knowledge workers waste up to 50% of their time dealing with mundane data quality issues and searching for reliable information [5]. These impacts can be categorized into several dimensions:

- *Financial Impact*: Beyond the direct costs of remediation, poor data quality leads to missed opportunities, regulatory fines, and inefficient operations. A study by IBM estimated that poor data quality costs the US economy over $3.1 trillion annually [6].
- *Decision Risk*: As organizations become more data-driven, the risk of making strategic decisions based on faulty data increases. This risk is particularly acute in domains like financial services, healthcare, and critical infrastructure.
- *Trust Erosion:* When data consumers encounter quality issues, they develop what some researchers call "data skepticism"—a tendency to question or even ignore data-driven insights, undermining the very culture organizations seek to build.
- *Operational Disruption:* Data pipeline failures can cascade through dependent systems, causing outages in critical business applications, delaying analytical insights, and consuming disproportionate engineering resources.

These consequences are intensified by what has been termed "data downtime"—periods when data is missing, inaccurate, or otherwise unfit for intended uses. According to Monte Carlo Data, organizations experience an average of 61 data incidents per month, with each incident taking an average of 13 hours to identify and resolve [7].

### C. Limitations of Traditional Approaches

Traditional approaches to data quality management have centered around three main strategies, each with significant limitations in modern data environments:

- *Manual Testing and Validation:* Many organizations rely on manual data quality checks performed by analysts or engineers. This approach fails to scale with data volume and velocity, covers only a fraction of potential issues, and shifts valuable human resources away from higher-value work.
- *Rule-Based Checks*: Static, rule-based validations (e.g., checking for nulls, duplicates, or out-of-range values) address only known failure modes and require constant maintenance as data schemas and semantics evolve. They also struggle with complex, cross-system dependencies.
- *Post-hoc Cleanup*: Reactive approaches that focus on cleaning data after issues have occurred fail to address root causes and create endless remediation cycles. They also typically occur too late to prevent downstream impacts.

These approaches share a fundamental limitation: they treat data quality as a static property rather than a dynamic system characteristic that requires continuous monitoring and adaptation. They also tend to focus on surface-level symptoms rather than underlying causes, leading to what some practitioners call "data quality whack-a-mole"—a never-ending cycle of fixing the same issues in different contexts.

## III. DATA OBSERVABILITY: FROM MONITORING TO UNDERSTANDING

### A. The Observability Paradigm

The concept of observability has its roots in control theory, where it refers to the ability to infer the internal state of a system based on its external outputs. In software engineering, observability emerged as an evolution of monitoring, moving beyond simple metrics to provide deeper insights into system behavior through the collection and analysis of logs, traces, and other telemetry data.

Data observability extends this paradigm to data pipelines and assets, providing visibility into the health, performance, and quality of data at each stage of its lifecycle. While monitoring answers the question "What is happening?", observability answers "Why is it happening?" This distinction is critical for building truly resilient data systems.

The five pillars of data observability, as articulated by Barr Moses and others [8], provide a framework for understanding this emerging discipline:

1) *Freshness:* Is data being updated at the expected cadence? Are there delays in processing or availability?

2) *Distribution:* Are the statistical properties of data (means, medians, cardinality, etc.) within expected ranges? Are there anomalies in the distribution that might indicate quality issues?

3) *Volume:* Is the amount of data being processed, stored, or queried within expected boundaries? Are there unexpected spikes or drops?

4) *Schema:* Are the structure, fields, and types of data consistent with expectations? Have there been unplanned or undocumented schema changes?

5) *Lineage:* How does data flow through systems? What are the dependencies between datasets, transformations, and consuming applications?

Together, these pillars enable data teams to build a comprehensive understanding of their data systems, identify potential issues before they impact downstream consumers, and trace problems to their root causes.

### B. From DevOps to DataOps

The emergence of data observability parallels the broader evolution from DevOps to DataOps—the application of DevOps principles and practices to data engineering and analytics. This evolution recognizes the unique characteristics of data workflows compared to traditional software development:

- *State Dependency:* Unlike stateless applications, data pipelines operate on and transform state. A single corruption can propagate through the entire system.

- *Semantic Complexity:* Understanding data requires domain knowledge beyond pure engineering expertise.
- *Asynchronous Feedback*: Issues in data pipelines may not manifest until long after they occur, making immediate feedback loops challenging.

DataOps practices emphasize continuous testing, monitoring, and improvement of data pipelines, creating a foundation for observability. However, true observability requires moving beyond process to implement technological solutions that provide the necessary visibility.

### C. Building Observable Data Systems

Implementing data observability requires a combination of architectural patterns, tooling, and organizational practices. Key elements include:

- *Instrumentation*: Adding metadata and telemetry collection points throughout data pipelines to capture operational metrics, data profiles, and quality indicators.
- *Centralized Observability Layer*: Creating a unified platform for collecting, storing, and analyzing observability data across the entire data ecosystem.
- *Semantic Layer*: Mapping technical metadata to business concepts to enable meaningful interpretation of quality metrics.
- *Anomaly Detection*: Implementing statistical and machine learning techniques to identify deviations from expected behavior without relying solely on predefined thresholds.

Organizations leading in this space, such as Spotify with its "Data Health System" and Netflix with its "Data Validation Framework," have demonstrated the value of investing in observability [9]. These companies report significant reductions in data incidents, faster time to resolution, and improved data team productivity. However, observability alone addresses only part of the challenge. To build truly resilient systems, organizations must move from observation to automated action—the domain of self-healing data pipelines.

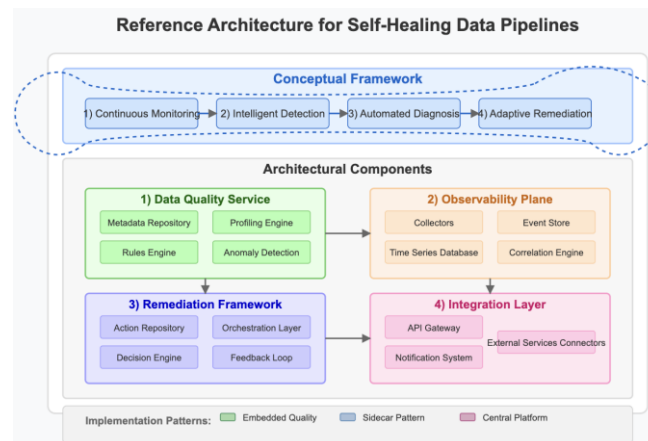## IV. REFERENCE ARCHITECTURE FOR SELF-HEALING DATA PIPELINES



*Fig1. Reference Architecture for Self-Healing Data Pipelines*

### A. Conceptual Framework

Self-healing data pipelines represent the convergence of data observability, quality automation, and intelligent remediation. The core principle is to create closed-loop systems that can detect, diagnose, and correct issues with minimal human intervention. This approach draws inspiration from concepts in control systems, site reliability engineering, and autonomous computing.

At a conceptual level, self-healing data pipelines incorporate four fundamental capabilities:

1) *Continuous Monitoring:* Real-time collection of metrics, profiles, and quality indicators across all data assets and pipeline components.

*2)* *Intelligent Detection:* Identification of anomalies, deviations, and potential issues through statistical analysis, machine learning, and domain-specific rules.

*3)* *Automated Diagnosis:* Determination of root causes through correlation analysis, impact assessment, and causal inference.

*4)* *Adaptive Remediation:* Implementation of corrective actions ranging from simple fixes (e.g., retrying failed jobs) to complex interventions (e.g., dynamically adjusting transformation logic).

These capabilities form a continuous loop, with each phase feeding into the next and creating a system that learns and improves over time. The degree of automation in the remediation phase can vary based on organizational risk tolerance, domain complexity, and technical maturity.

*B. Architectural Components*

Translating this conceptual framework into a practical architecture requires several key components:

*1)* *Data Quality Service*

The data quality service serves as the central nervous system of the self-healing architecture. It combines:

• *Metadata Repository:* A comprehensive catalog of data assets, their schemas, relationships, and expected quality characteristics.

• *Rules Engine:* A system for defining, managing, and executing data quality rules, both predefined and dynamically generated.

• *Profiling Engine:* Automated analysis of data characteristics, distributions, and patterns to establish baselines and detect anomalies.

• *Anomaly Detection:* Machine learning models that identify unexpected deviations from historical patterns and expected behaviors.

Modern implementations often leverage graph databases for metadata management, streaming technologies for real-time profiling, and ensemble machine learning approaches for anomaly detection.

*2)* *Observability Plane*

The observability plane collects, stores, and processes telemetry data from across the data ecosystem:

• *Collectors*: Agents and connectors that gather metrics, logs, and events from data sources, processing engines, and pipeline components.

• *Time Series Database*: Optimized storage for high-volume metric data with efficient query capabilities.

• *Event Store*: Persistent storage for logs, alerts, and significant events with support for complex queries.

• *Correlation Engine*: Analysis tools for identifying relationships between events, metrics, and quality indicators.

Leading organizations have found that the volume of observability data often exceeds that of the business data itself, requiring specialized infrastructure and retention policies.

*3)* *Remediation Framework*

The remediation framework enables automated response to detected issues:

• *Action Repository*: A catalog of predefined remediation actions with their preconditions, effects, and implementation details.

• *Decision Engine:* Logic for selecting appropriate remediation actions based on issue type, severity, and context.

• *Orchestration Layer*: Execution environment for implementing remediation actions across distributed systems.

• *Feedback Loop*: Mechanisms for tracking the effectiveness of remediation actions and refining future responses.

The sophistication of remediation actions can range from simple "retry logic" to complex interventions involving dynamic reconfiguration of pipeline components or even synthesis of corrective transformations.

*4) Integration Layer*

The integration layer connects the self-healing system with the broader data ecosystem:

- *API Gateway*: Standardized interfaces for interacting with the self-healing system from external tools and applications.
- *Notification System*: Channels for alerting stakeholders about detected issues, implemented remediation, and required human intervention.
- *External Services Connectors*: Integration points with related systems such as data catalogs, workflow orchestrators, and monitoring platforms.

This layer enables the self-healing system to function as part of a cohesive data management ecosystem rather than an isolated solution.

*C. Implementation Patterns*

While the reference architecture provides a blueprint, implementation approaches vary based on organizational context, existing infrastructure, and maturity level. Three common patterns have emerged:

- *Embedded Quality*: Integrating observability and remediation capabilities directly into data processing frameworks. For example, modern frameworks like Apache Spark and dbt include built-in functionality for data quality validation and error handling.
- *Sidecar Pattern*: Deploying observability and remediation components alongside existing pipeline components without modifying them. This pattern is particularly valuable for legacy systems or third-party components that cannot be directly modified.
- *Central Platform*: Implementing a dedicated data quality platform that connects to all data assets and pipelines through standardized interfaces. This approach provides the most comprehensive coverage but requires significant investment in integration.

Organizations often combine these patterns, applying different approaches to different parts of their data ecosystem based on criticality, complexity, and technical constraints.

Use the enter key to start a new paragraph. The appropriate spacing and indent are automatically applied.

## V. IMPLEMENTATION STRATEGIES AND CASE STUDIES

*A. Phased Implementation Approach*

Building self-healing data pipelines represents a significant transformation for most organizations. A phased implementation approach helps manage complexity and demonstrate value incrementally:

*1) Foundation Phase:* Establish basic observability by instrumenting critical data assets and pipelines. Focus on collecting metadata, profiling data, and establishing baselines for normal behavior.

*2) Detection Phase:* Implement anomaly detection and quality validation rules. Begin generating alerts for potential issues, even if remediation remains manual.

*3) Diagnosis Phase:* Develop correlation capabilities to identify root causes and impact assessments for detected issues. Implement automated root cause analysis for common failure patterns.

*4) Controlled Remediation:* Implement automated remediation for low-risk, well-understood issues. Establish clear boundaries between automated and human-in-the-loop scenarios.

*5) Advanced Automation:* Gradually expand the scope of automated remediation, incorporating machine learning for adaptive responses and handling increasingly complex scenarios.

This approach allows organizations to build capability and confidence progressively, learning and adjusting as they go. It also helps data teams adapt to new workflows and responsibilities.

*B. Case Study: Financial Services Firm*

A global financial services firm implemented self-healing data pipelines to address recurring quality issues in their customer data integration processes. Their implementation journey illustrates several key success factors:

*1) Initial Challenge*

The firm processed data from over 200 source systems, integrating it into a central customer data platform that powered regulatory reporting, analytics, and customer-facing applications. They experienced frequent quality issues, including:

- Missing or delayed data from critical sources
- Schema drift causing transformation failures
- Duplication and inconsistency across sources
- Reference data integrity problems

These issues resulted in approximately 350 hours of data downtime per month, required 12 full-time engineers for troubleshooting, and significantly delayed regulatory reporting.

*2) Implementation Approach*

The firm adopted a phased approach focused initially on their most critical data domains:

1. They deployed data observability tools to monitor data freshness, volume, and schema stability across all ingestion points.
2. They implemented anomaly detection for key quality metrics, generating alerts when data deviated from expected patterns.
3. They developed a remediation framework with automated actions for common issues:
   o Automatic retries for failed connections with exponential backoff
   o Dynamic schema adjustment for non-breaking changes
   o Automated deduplication based on configurable business rules
   o Source priority rules for resolving conflicts in overlapping data
4. They established a feedback loop, tracking the effectiveness of automated remediation and continuously refining detection and correction logic.

*3) Results*

After 12 months, the firm reported:

- 82% reduction in data downtime
- 76% decrease in manual intervention for data quality issues
- 15-hour improvement in regulatory reporting timeliness
- Redeployment of 8 engineers from firefighting to strategic initiatives

The team identified three critical success factors:

1. Starting with a focused scope rather than attempting to implement self-healing across all systems simultaneously
2. Establishing clear thresholds for automated vs. human-in-the-loop remediation
3. Creating a continuous learning process to refine detection and remediation logic

*C. Case Study: Healthcare Analytics Provider*

A healthcare analytics provider implemented self-healing pipelines to ensure reliability in their clinical data processing workflows. Their experience highlights different aspects of the implementation journey:

*1)    Initial Challenge*

The provider ingested clinical data from hundreds of healthcare organizations, standardized it to common formats, and produced analytics models for quality improvement, population health, and clinical research. They faced unique challenges:

- Extreme heterogeneity in source data formats and quality
- Strict regulatory requirements for data handling and privacy
- Complex transformation logic requiring clinical domain knowledge
- High stakes for data accuracy in clinical decision support

Manual quality management processes were overwhelming their team and creating unacceptable delays in data availability.

*2)    Implementation Approach*

The provider took a domain-centric approach to implementation:

1.     They created a domain-specific quality framework incorporating clinical knowledge, defining expected relationships between medical concepts, valid value ranges, and consistency rules.
2.     They implemented a metadata-rich observability layer that tracked not just technical metrics but semantic characteristics of the data.
3.     They developed a tiered remediation approach:
o          Tier 1: Fully automated remediation for well-understood technical issues
o          Tier 2: Suggested remediation with human approval for domain-specific issues
o          Tier 3: Alert-only for complex issues requiring clinical expertise
4.     They implemented a learning system that captured expert decisions for Tier 3 issues, gradually building a knowledge base that enabled more issues to be handled at Tier 2 or Tier 1.

*3)    Results*

After 18 months, the provider reported:

- 94% of technical quality issues resolved automatically
- 67% of domain-specific issues handled through suggested remediation
- 3x improvement in time-to-value for new data sources
- Expanded coverage to 3x more healthcare organizations without increasing quality management staff

Their key insights included:

1.     The importance of domain-specific quality rules beyond generic data quality dimensions
2.     The value of progressive automation that learns from expert intervention
3.     The need to balance automation with appropriate human oversight in high-stakes domains

## VI. FUTURE DIRECTIONS AND EMERGING TECHNOLOGIES

### A. Machine Learning for Adaptive Quality Management

The next frontier in self-healing data pipelines involves more sophisticated applications of machine learning beyond basic anomaly detection:

- *Predictive Quality Management*: Using leading indicators to predict potential quality issues before they manifest, enabling truly preventative actions.
- *Automatic Rule Generation*: Learning quality rules from data rather than requiring explicit definition, helping systems adapt to evolving data characteristics.
- *Causal Inference*: Moving beyond correlation to establish causal relationships between system events and data quality issues, enabling more precise diagnosis.
- *Semantic Understanding*: Developing models that understand the business meaning of data, enabling quality assessment based on fitness for specific use cases rather than generic technical criteria.

Research in these areas is advancing rapidly, with organizations like MIT's Data Systems and AI Lab (DSAIL) and companies like Anomalo and Monte Carlo leading innovation [10].

*B. Federated Data Quality*

As data ecosystems become increasingly distributed—spanning multiple clouds, on-premises systems, and edge environments—a new approach to data quality is emerging: federated data quality management. This approach:

- Distributes quality verification and enforcement to where data resides
- Establishes consistent quality standards across heterogeneous environments
- Aggregates quality signals across the distributed landscape
- Enables coordinated remediation across system boundaries

Early implementations of federated approaches show promise in complex, multi-cloud environments where centralized quality management would create unacceptable latency or governance challenges.

*C. Quality-Native Data Systems*

Looking further ahead, we anticipate the emergence of "quality-native" data systems—platforms that treat quality as a first-class concern rather than an afterthought. These systems will:

- Embed quality verification into core data operations
- Provide "quality guarantees" similar to the consistency guarantees in distributed databases
- Automatically adjust processing based on quality characteristics
- Support quality-based routing and prioritization of data flows

Just as "cloud-native" technologies reimagined applications for distributed environments, quality-native systems will fundamentally rethink data processing around the central importance of quality and reliability.

## VII. CONCLUSION AND RECOMMENDATIONS

*A. The Path to Data Reliability*

Data observability and self-healing pipelines represent a fundamental shift in how organizations manage data quality—from reactive remediation to proactive prevention, from manual intervention to intelligent automation. This shift is not merely technical but organizational, requiring new skills, processes, and mindsets.

The journey toward self-healing data pipelines is, in many ways, a journey toward what some have called "Data Reliability Engineering" (DRE)—the application of site reliability engineering principles to data systems [11]. Just as SRE transformed application management, DRE promises to transform data management by establishing clear reliability objectives, implementing automated solutions, and continuously improving both systems and processes.

*B. Key Recommendations*

Based on the research and case studies presented in this paper, we offer several recommendations for organizations embarking on this journey:

*1) Start with Observability:* Build comprehensive visibility into your data ecosystem before attempting automated remediation. You cannot fix what you cannot see.

*2) Adopt Domain-Driven Quality:* Define quality not just in technical terms but in relation to business domains and use cases. The same data may have different quality requirements depending on how it's used.

*3) Embrace Progressive Automation:* Begin with alert-only systems, then graduate to suggested remediation, and finally to fully automated healing where appropriate. Build confidence in your automation incrementally.

*4) Invest in Metadata:* Rich, accurate, and comprehensive metadata is the foundation for effective observability and remediation. Prioritize tools and processes that capture and maintain high-quality metadata.

*5)* *Balance Technical and Organizational Change:* Implementing self-healing pipelines requires both technological solutions and new organizational practices. Invest in both dimensions equally.

*6)* *Measure Impact Holistically:* Evaluate success not just through technical metrics like downtime reduction but through business impacts like improved decision-making, accelerated insights, and enhanced trust in data.

The organizations that successfully implement these recommendations will not only reduce costs and improve efficiency but also unlock new possibilities for data-driven innovation. By ensuring that data is consistently reliable, timely, and accurate, they create the foundation for advanced analytics, machine learning, and artificial intelligence initiatives that depend on high-quality data.

As data continues to grow in volume, velocity, and importance, the ability to automatically ensure its quality and reliability will become not just a competitive advantage but a fundamental requirement for organizational success. The path to self-healing data pipelines represents one of the most important journeys that data-driven organizations can undertake.

**REFERENCES**

[1] T. C. Redman, "Data's Credibility Problem," *Harvard Business Review*, Dec. 2013. https://hbr.org/2013/12/datas-credibility-problem

[2] Harvard Business Review, "Getting Serious About Data and Data Science," Oct. 2020. https://hbr.org/2020/10/getting-serious-about-data-and-data-science

[3] L. Goasduff, "Data and Analytics Leaders Face Acute Data Quality Issues," *Gartner*, May 2022. https://www.gartner.com/en/articles/data-and-analytics-leaders-face-acute-data-quality-issues

[4] K. Chui, M. Manyika, and M. Miremadi, "What AI can and can't do (yet) for your business," *McKinsey Quarterly*, Jan. 2018. https://www.mckinsey.com/business-functions/mckinsey-analytics/our-insights/what-ai-can-and-cant-do-yet-for-your-business

[5] T. Davenport and R. Bean, "Big Companies Are Embracing Analytics, But Most Still Don't Have a Data-Driven Culture," *Harvard Business Review*, Feb. 2018.https://hbr.org/2018/02/big-companies-are-embracing-analytics-but-most-still-dont-have-a-data-driven-culture

[6] IBM Cloud Education, "Big Data Analytics," *IBM*, June 2020. https://www.ibm.com/cloud/learn/big-data-analytics

[7] AWS, "Data Quality in Machine Learning," *AWS Documentation*, 2023.https://docs.aws.amazon.com/wellarchitected/latest/machine-learning-lens/data-quality.html

[8] Microsoft, "Data quality and preparation," *Microsoft Learn*, 2023. https://learn.microsoft.com/en-us/azure/architecture/data-guide/technology-choices/data-quality

[9] F. Olumuyiwa, "Data Quality at Netflix: Detecting Data Anomalies at Scale," *Netflix Technology Blog*, Nov. 2019. https://netflixtechblog.com/data-quality-at-netflix-detecting-data-anomalies-at-scale-41cbf081b1b

[10] AWS, "Building a serverless data quality framework," *AWS Big Data Blog*, July 2021. https://aws.amazon.com/blogs/big-data/building-a-serverless-data-quality-framework-using-aws-glue-and-aws-lambda/

[11] B. Beyer et al., "Site Reliability Engineering," *Google SRE*, 2016. https://sre.google/sre-book/table-of-contents/