



Mission success is only as reliable as your data

The essential guide to data reliability in the AI era.

Table of contents

Fragmented governance		Six steps to data reliability	9
is your biggest obstacle	3	Step 1: Profile the data	1C
Why unified governance		Step 2: Define data policies and business rules	10
is the answer	3	Step 3: Detect anomalies	11
What is Data Confidence™?	4	Step 4: Monitor for impact	11
Getting prepared		Step 5: Notify key experts	12
for data reliability	5	Step 6: Optimize continuously	12
A solution to complex data challenges	7	Collibra: Feel confident about your data and your Al	13



Fragmented governance is your biggest obstacle

The gap between what federal agencies need to accomplish with data and Al and what they're actually capable of is widening. Further widening that gap are the wasted resources caused by unreliable data. With 68% of federal Chief Data Officers (CDOs) saying they have less than 10 FTEs to fulfill their mission, it's no wonder why it's difficult to identify data anomalies and arm agency employees with the data needed to reduce waste and enhance citizen services.1

In most agencies, data exists in isolated pockets across apps. multiple clouds and on-premises systems. This creates blind spots in what data exists, who has access and how it's being used. Traditional approaches that tether governance to specific systems or platforms won't cut it anymore—especially as Al multiplies the risks of unreliable and noncompliant use.

The answer isn't another point solution—it's unified governance across your agency's entire ecosystem.

Why unified governance is the answer

The agencies that can unify governance across every data source, use case and user will be ready to act on the best Al opportunities first—while controlling Al risks before it's too late. That's why federal agencies are turning to unified governance that gives them visibility, context and control throughout the full data cycle.

Almost half of federal **CDOs say they need** funding / staffing to achieve their mission²

¹ https://www2.deloitte.com/content/dam/Deloitte/us/Documents/public-sector/2024-cdo-survey.pdf 2 https://www2.deloitte.com/content/dam/Deloitte/us/Documents/public-sector/2024-cdo-survey.pdf



What is Data **Confidence?**

Data Confidence is the way you and your colleagues feel when your agency can accelerate every data and Al use case that supports your mission — without compromising on efficiency, safety or quality.

It happens when governance becomes an enabler rather than a bottleneck. Your people can find, understand and use trusted data across every system. Business context flows alongside technical metadata. And policies apply consistently everywhere data lives.

Bottom line: When your people can trust, comply and consume data confidently, innovation accelerates. That's Data Confidence.

Unreliable data can be costly to your organization, including:

- Inaccurate risk evaluations: Using the wrong data can mean the cost of improper balancing of financial, operational and security issues
- Unsustainable resource allocations: Unreliable data can mean spending too much on budgeting, staffing or inventory. It can also mean missed opportunities
- **Erroneous performance evaluations:** Missing or unreliable data can result in a lack of accountability, losses to competition and outcomes that were never in line with stakeholders expectations in the first place

Getting prepared for data reliability

You can improve the reliability of your data with an effective data quality solution. To get started, you'll want to ask yourself and your colleagues a few fundamental questions.

Question 1.

Are data quality and reliability "events" understood when they occur? How commonplace are they?

Understanding the frequency and nature of data quality issues is crucial, especially as Al investments expand and regulatory landscapes evolve. This insight helps determine the effectiveness of your current data management practices and highlights areas needing attention.

Question 2. What data quality dimensions or metrics show success for you?

Identify the key metrics that reflect data quality in your organization. This could include accuracy, completeness, consistency, timeliness and validity. Knowing these will help you set benchmarks and track improvements over time, ensuring your data meets the standards required for Al and compliance purposes.

Question 3.

Who across your organization must have visibility on the reliability of data?

It's essential to identify the stakeholders who need access to data quality insights. This could include data engineers, analysts, business leaders and compliance officers. Ensuring that the right people have visibility into data reliability helps in making informed decisions and maintaining accountability, which is increasingly important in a regulated environment.

Question 4.

Do those involved with business and rule validation have technical backgrounds?

Understanding the technical expertise of your team members involved in data governance is crucial. This will help in designing training programs and support systems to bridge any knowledge gaps and ensure effective rule validation, which is vital as Al integration becomes more complex.

Question 5.

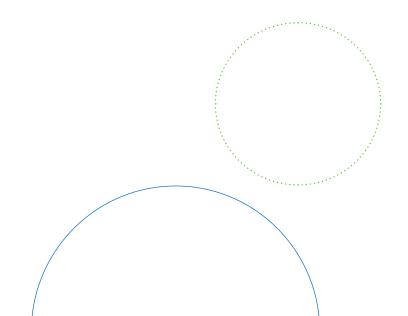
How manual is the process of generating data quality rules today?

Assess the current level of automation in your data quality management processes. Manual processes can be error-prone and time-consuming. Identifying opportunities for automation can enhance efficiency and accuracy, which is essential for keeping pace with the rapid changes in data volume and regulatory requirements.

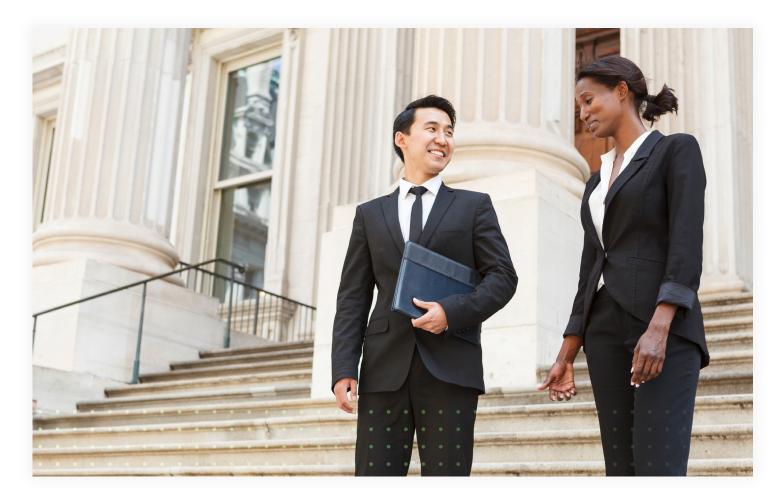
These questions will guide you in evaluating your current data quality and observability readiness. And, they'll help you pinpoint where you need focus and improvements, setting a solid foundation for robust data governance and Alimplementation.

Why Collibra Data Quality

- Monitor data quality issues across many sources with automated and targeted rule writing
- Remediate faster with no-code, business ready DQ dimensions and self-service rules for accurate and consistent data products
- Visualize data health and certify data for trusted business decisions



A solution to complex data challenges



Traditional manual processes are no longer sufficient to manage the vast amounts of data and the rapid pace of change. This is where machine learning (ML) can play a transformative role. By leveraging ML, organizations can automate and enhance their data observability practices, ensuring reliable data quality, regulatory compliance and improved decision-making.

However, it's important to take a structured approach when implementing machine learning to tackle data observability challenges.

Here are the six key steps to ensure your team's success:

- 1. Profile the data: Discover and classify your data sources, focusing on sensitive data types
- Define data policies and business rules: Implement continuous testing and validation mechanisms to maintain data integrity

- 3. Detect anomalies: Use machine learning to monitor data and identify deviations from normal behavior
- 4. Monitor for impact: Correlate anomalies with changes and events to understand their impact
- 5. Notify key experts: Alert relevant stakeholders to initiate remediation processes
- 6. Optimize continuously: Evolve policies, rules and reports to enhance data quality goals

Ready to move forward on your journey to data reliability? In the pages to follow, we'll expand on each of these steps, identify common obstacles and provide four questions you can ask to clear your path to data reliability.

The six steps to data reliability



Step 1: Profile the data

Start by discovering all your data sources and classifying them, with a particular focus on identifying sensitive data types. Profiling your data helps in understanding the nature, structure, and quality of your data assets. This foundational step is crucial for setting the stage for effective data governance and observability, especially as you increase your Al investments and the regulatory landscape evolves.

Common obstacles

- Data silos: Data spread across different departments or systems can be challenging to consolidate
- Data volume: Large volumes of data can make profiling and classification time-consuming and resource-intensive
- Data complexity: Different data formats and structures add to the complexity of profiling

Four questions to ask

- Where is our data stored, and how can we access it?
- What types of data do we have, and which are sensitive or regulated?
- How complete and accurate is our data? 3.
- What tools and processes do we need to efficiently profile our data?

Step 2: Define data policies and business rules

Based on the types of data you have, define policies and business rules that guide how data should be handled. Implement continuous testing and validation mechanisms to ensure data pipelines do not contain data that violates these policies. This proactive approach helps your agency maintain data integrity, delivering better decision-making for mission success.

Common obstacles

- Lack of standardization: Inconsistent data policies across the organization
- Regulatory complexity: Navigating the myriad of evolving data regulations
- **Resistance to change:** Difficulty in getting buy-in from all stakeholders

Four questions to ask

- What are the critical data policies and business rules needed for our data types?
- 2. How can we ensure continuous compliance with these policies?
- 3. What are the potential risks if these rules are violated?
- 4. Who is responsible for enforcing these policies, and how do we ensure accountability?



Step 3: Detect anomalies

Establish a baseline for normal behavior within your data. Use machine learning algorithms to monitor data continuously and detect deviations from this baseline. Anomaly detection is vital for identifying potential data quality issues before they escalate into bigger problems, which is increasingly important as Al applications depend on high-quality data to function correctly and comply with regulations.

Common obstacles

- Defining normal behavior: Establishing accurate baselines can be complex
- Algorithm selection: Choosing the right ML algorithms for effective anomaly detection
- False positives: High rates of false positives can overwhelm teams and reduce trust in the system

Four questions to ask

- What constitutes normal behavior for our data?
- Which ML algorithms are best suited for our anomaly detection needs?
- How do we balance sensitivity and specificity to minimize false positives?
- How often should we review and update our baseline metrics?

Step 4: Monitor for impact

Once anomalies are detected, correlate these deviations with unintended changes and other events to identify the root cause and assess the potential impact. Understanding the impact of data quality issues on your operations helps in prioritizing remediation efforts effectively, ensuring your Al models and analytics applications continue to perform accurately and reliably.

Common obstacles

- Correlation complexity: Difficulties in accurately correlating anomalies with their impacts
- Resource allocation: Limited resources to investigate and resolve detected issues
- **Impact assessment:** Challenges in quantifying the business impact of data quality issues

Four questions to ask

- How do we identify and document the root causes of anomalies?
- 2. What tools can help us correlate anomalies with their impacts efficiently?
- 3. How do we prioritize issues based on their potential impact?
- What metrics should we use to assess the impact of data quality issues?

Step 5: Notify key experts

Provide contextual alerts to relevant stakeholders—including data engineers, analysts, division leaders and compliance officers—to initiate remediation processes. Timely and informed notifications ensure that the right people can take swift action to resolve data quality issues, maintaining both operational efficiency and regulatory compliance.

Common obstacles

- Alert fatigue: Too many alerts can lead to important notifications being ignored
- Contextual information: Ensuring alerts contain enough context for quick action
- Stakeholder engagement: Ensuring all relevant stakeholders are engaged and responsive

Four questions to ask

- Who needs to be notified when a data quality issue is detected?
- 2. What information should be included in the alerts to make them actionable?
- How can we ensure timely responses to critical alerts?
- What is our process for escalating unresolved issues?

Step 6: Optimize continuously

Evolve your data policies, rules, and reports based on insights gained from monitoring and anomaly detection. Continuous optimization ensures that your data governance practices remain effective and aligned with evolving data quality goals, supporting the scalability of Al initiatives and adherence to changing regulatory standards.

Common obstacles

- **Continuous improvement:** Keeping up with the need for ongoing adjustments and optimizations
- Feedback loop: Establishing effective feedback mechanisms to inform improvements
- Regulatory changes: Adapting to new and changing regulations in a timely manner

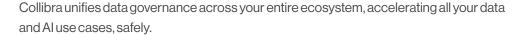
Four questions to ask

- How do we gather and analyze feedback to improve our data policies and rules?
- 2. What process do we have in place for updating our data governance practices?
- 3. How do we stay informed about regulatory changes and ensure compliance?
- 4. How can we measure the effectiveness of our optimizations over time?

The path to data reliability is here

By following these steps, your organization can leverage machine learning to enhance data observability, ensuring high data quality and compliance. This structured approach leads to better decision-making, reduced costs associated with poor data quality, and improved overall data management.

Collibra: Feel confident about your data and your Al



We automate the tedious tasks so your team can focus on what matters. Bringing technical and business users together, with natural language rule writing and AI-powered assistance that makes trusted, high-quality data accessible to everyone.

Most importantly, we help you build Data Confidence—that state where everyone in your organization can trust, comply and consume data without fear.

Is your organization ready for Data Confidence? Discover Collibra.

Collibra Data Quality & Observability (DQ&O) leverages Adaptive Rules for intelligent monitoring, providing complete visibility into technical metrics like null checks, row counts and outliers. Our machine learning algorithms ensure these rules are always up-to-date by learning from past and new data. This comprehensive approach connects data issues to their root causes, linking data ownership, lineage and detailed reliability analysis within our platform.

We offer direct integration for data quality across the data catalog, a unique feature among vendors. Our anomaly detection capabilities provide thorough outlier detection, and our health reporting tools offer business leaders like you clear insights into data quality dimensions specific to your needs, without overwhelming you and your data professionals with information.

Collibra's DQ Pushdown enhances efficiency, cost management, security and reduces time to value. Our break record storage and quarantine options for exception records ensure safe and secure handling of data remediation issues.

Learn more about Collibra Data Quality & Observability.



Ready to ensure data reliability in your organization? Discover how Collibra can transform your data quality and observability practices. Learn more about Collibra Data Quality & Observability.