
| RESEARCH ARTICLE

Leveraging AI for Better Data Quality and Insights

Ankit Pathak

Indian Institute of Technology (Indian School of Mines), India

Corresponding Author: Ankit Pathak, **E-mail:** contact.ankitpathak@gmail.com

| ABSTRACT

The exponential growth of data across industries has highlighted the critical importance of data quality management for ensuring reliable insights and decision-making. Artificial intelligence has emerged as a transformative force in this domain, offering sophisticated approaches to detect errors, inconsistencies, and anomalies in complex datasets. This article explores the fundamental principles of data quality control, examines AI-powered methodologies including machine learning algorithms, deep learning architectures, and natural language processing techniques, and investigates their domain-specific applications across healthcare, finance, marketing, manufacturing, and government sectors. Despite significant advancements, challenges persist related to scalability, human-AI collaboration, privacy concerns, model interpretability, and adaptation to evolving data patterns. Emerging trends such as explainable AI, human-in-the-loop frameworks, transfer learning, federated approaches, real-time monitoring, and quantum computing applications promise to further enhance AI's effectiveness in elevating data quality standards and unlocking greater value from organizational data assets.

| KEYWORDS

Data quality dimensions, Anomaly detection, Natural language processing, Entity resolution, Privacy-preserving techniques

| ARTICLE INFORMATION

ACCEPTED: 09 April 2025

PUBLISHED: 03 May 2025

DOI: 10.32996/jcsts.2025.7.3.33

1. Introduction

In today's data-driven landscape, organizations are generating and collecting unprecedented volumes of information. According to recent analysis, the global datasphere is projected to grow from 33 zettabytes in 2018 to 175 zettabytes by 2025, with enterprises creating and managing 60% of the world's data [1]. This exponential growth has transformed how businesses operate across sectors, with an increasing number of organizations now depending on data-driven insights for strategic decision-making.

The utility of these massive datasets fundamentally depends on their quality. Research indicates that poor data quality costs organizations millions annually, affecting everything from operational efficiency to customer satisfaction. As data volumes continue to expand, maintaining high-quality standards becomes increasingly challenging yet essential for deriving meaningful insights and competitive advantage.

Artificial intelligence has emerged as a transformative force in addressing complex data quality challenges. Traditional approaches struggle with the volume, variety, velocity, and veracity characteristics of modern big data environments [2]. These "4 Vs" represent significant hurdles that conventional data management systems cannot effectively overcome without intelligent automation. AI techniques offer adaptable and scalable alternatives by leveraging machine learning algorithms and neural networks to detect patterns and anomalies that would be impossible to identify through manual inspection.

The application of AI to data quality has shown promising results across various domains. Predictive models can now identify potential issues before they impact downstream processes, while automated cleansing systems significantly reduce manual

preparation efforts. Furthermore, AI-based entity resolution techniques have achieved high match accuracy rates in complex datasets with millions of records, outperforming traditional approaches.

This article aims to provide a comprehensive examination of the intersection between artificial intelligence and data quality management. Our objectives include establishing a foundational understanding of data quality principles, analyzing various AI techniques currently employed to enhance data quality, and exploring emerging trends in this rapidly evolving field. We will examine both theoretical frameworks and practical implementations, drawing from academic research and industry experience.

Our analysis encompasses the full spectrum of AI technologies and their specific applications in addressing various dimensions of data quality, including accuracy, completeness, consistency, and timeliness. Additionally, we investigate the organizational and technical challenges associated with implementing AI-powered data quality solutions. By providing this perspective, we aim to equip readers with knowledge necessary to leverage AI effectively for improving data quality and driving more reliable decision-making.

2. Fundamental Principles of Data Quality Control

2.1. Data Quality Dimensions: Accuracy, Completeness, Consistency, and Timeliness

Data quality is a multidimensional concept encompassing several key attributes that collectively determine the fitness of data for its intended use. According to Myles Suer's comprehensive analysis, organizations implementing structured data quality programs report up to a 40% reduction in time spent on data preparation tasks [3]. The primary dimensions include accuracy (the degree to which data correctly represents real-world entities), completeness (the presence of all required data points), consistency (the absence of contradictions across datasets), and timeliness (the availability of data when needed). These dimensions serve as the foundation for establishing measurable quality metrics and improvement targets.

Additional dimensions gaining recognition include relevance (alignment with business needs), uniqueness (absence of duplication), validity (conformance to syntax rules), and accessibility (ease of obtaining data). Each dimension requires specific measurement techniques and improvement strategies, with the relative importance varying based on organizational context and specific use cases.

2.2. Data Profiling: Understanding Your Data Landscape

Data profiling forms the foundation of effective quality management by providing comprehensive visibility into the structural characteristics, content, and relationships within datasets. This process involves analyzing data to discover patterns, identify anomalies, and understand distributions—essential insights that guide subsequent quality improvement efforts. Recent research published in the *Journal of Big Data* indicates that organizations implementing systematic profiling practices identify and address quality issues 37% faster than those using ad hoc approaches [4].

Modern data profiling tools employ statistical analysis, pattern recognition, and relationship discovery techniques to generate metadata about column properties, cross-column dependencies, and business rule compliance. Through profiling, organizations can establish quality baselines, prioritize remediation efforts, and measure improvement over time. The insights gained from profiling also inform data model refinements, cleansing rules, and validation criteria for ongoing quality assurance.

2.3. Data Cleansing: Techniques for Error Detection and Correction

Data cleansing encompasses methodologies for identifying and addressing quality defects across structured and unstructured data assets. This process targets various error types, including missing values, duplicates, outliers, inconsistent formats, and rule violations. As highlighted by Myles Suer, organizations with mature data cleansing practices report a 65% increase in user trust in their data assets, leading to more confident decision-making and analytics adoption [3].

Modern cleansing approaches combine automated techniques with domain expertise for optimal results. Common methodologies include standardization, enrichment, deduplication, and imputation. While traditional rule-based cleansing remains prevalent, machine learning algorithms increasingly augment these approaches by detecting complex patterns and suggesting corrections based on historical data.

2.4. Data Validation: Ensuring Compliance with Business Rules and Standards

Data validation ensures that information meets predefined quality criteria before entering operational systems or analytical environments. This process verifies compliance with syntactic rules and semantic constraints. According to the *Journal of Big Data*, organizations implementing automated validation frameworks experience a 44% reduction in the time required to identify and resolve data quality issues [4].

Effective validation strategies operate at multiple levels, including field-level checks, record-level rules, cross-record constraints, and dataset-level validations. Modern validation frameworks increasingly incorporate machine learning to develop adaptive rules

that evolve with changing data patterns and business requirements, moving beyond static, predefined checks toward more dynamic, context-aware validation.

2.5. The Data Quality Management Lifecycle

The data quality management lifecycle provides a structured framework for continuous quality improvement across organizational data assets. This cyclical process begins with defining quality requirements based on business objectives, followed by assessment, remediation, prevention, and monitoring phases. Myles Suer's research indicates that organizations with formalized lifecycle management achieve a 30% higher return on their data investments compared to those with fragmented approaches [3].

The lifecycle typically encompasses requirements gathering, assessment, remediation, prevention, and monitoring. Effective implementation requires clear governance structures, defined roles and responsibilities, appropriate technology enablement, and cultural commitment to data quality as an organizational priority.

Data Quality Practice	Performance Impact	Improvement Percentage
Structured Data Quality Programs	Reduction in Data Preparation Time	40%
Systematic Profiling Practices	Faster Identification of Quality Issues	37%
Mature Data Cleansing Practices	Increase in User Trust	65%
Automated Validation Frameworks	Reduction in Resolution Time	44%
Formalized Lifecycle Management	Higher Return on Data Investments	30%

Table 1: Impact of Data Quality Management Practices on Organizational Performance [2, 4]

3. AI-Powered Approaches to Data Quality Control

3.1. Machine Learning for Anomaly Detection

Machine learning techniques have revolutionized anomaly detection in data quality management, offering powerful methods to identify outliers, inconsistencies, and patterns that deviate from expected norms. According to recent research in the Journal of Manufacturing Systems, organizations implementing ML-based anomaly detection experience a significant reduction in false alarms compared to traditional statistical process control methods, with a notable 30% decrease in manufacturing environments [5]. These techniques are particularly valuable when dealing with high-dimensional data where conventional statistical methods struggle to capture complex relationships.

Supervised Learning Approaches

Supervised learning algorithms leverage labeled training data to learn patterns distinguishing normal from anomalous records. Classification techniques such as Random Forests and Support Vector Machines have demonstrated high effectiveness in identifying manufacturing defects and data quality issues when provided with sufficient high-quality training examples. These approaches excel when historical data with known quality issues is available, enabling the model to learn from past errors and apply these lessons to new data.

Unsupervised Learning for Pattern Recognition

Unsupervised learning techniques identify anomalies without requiring labeled training data, making them particularly valuable in environments where quality issues are not previously documented or are constantly evolving. Recent advancements in manufacturing fault detection using unsupervised learning show promising results in detecting production line irregularities with minimal prior knowledge [5].

Semi-Supervised Learning Techniques

Semi-supervised approaches combine elements of both supervised and unsupervised learning, utilizing a small set of labeled data alongside a larger unlabeled dataset. This hybrid approach addresses the common challenge of limited labeled examples while leveraging the patterns present in abundant unlabeled data, proving especially valuable in manufacturing environments with sparse historical defect data.

3.2. Deep Learning Architectures for Complex Data Quality Problems

Deep learning architectures have emerged as powerful tools for addressing complex data quality challenges that resist solutions through traditional methods. Their ability to automatically extract hierarchical features and model intricate relationships makes them particularly suited to quality issues in unstructured or semi-structured data.

Neural Networks for Missing Value Imputation

Neural network architectures have demonstrated remarkable capabilities in accurately imputing missing values across diverse data types. In sensor networks, where data quality is critical for accurate monitoring and control, neural imputation methods have shown superior performance compared to traditional techniques, especially in scenarios with multiple missing values [6].

Autoencoders for Data Denoising

Autoencoders excel at noise reduction and data reconstruction by learning compressed representations of data. These models identify and preserve essential patterns while filtering out noise or inconsistencies. In IoT sensor networks, autoencoder-based approaches have proven particularly effective for addressing noise contamination and inconsistencies in environmental monitoring data [6].

Recurrent Neural Networks for Temporal Data Quality

Recurrent Neural Networks, particularly LSTM variants, have proven highly effective for quality control in temporal or sequential data. Research in intelligent sensor systems demonstrates that LSTM models can effectively capture time-series dependencies to detect anomalies and improve data reliability in environmental monitoring applications [6].

3.3. Natural Language Processing for Textual Data Quality

Natural Language Processing techniques have transformed quality management for textual data, enabling organizations to systematically assess and improve unstructured information.

Entity Recognition and Resolution

Named Entity Recognition systems identify and categorize key elements in text, facilitating standardization and quality control. Advanced manufacturing systems increasingly use NLP-based approaches to extract critical information from maintenance logs and production documentation, improving data standardization across operations [5].

Text Normalization and Standardization

Text normalization processes transform textual data into consistent formats, addressing variations in spelling, abbreviations, and syntactic structures that complicate analysis and integration. In manufacturing contexts, these techniques help standardize documentation across multiple facilities and supply chain partners.

Semantic Consistency Checking

Semantic analysis techniques evaluate the meaning and contextual appropriateness of textual data, enabling quality assessment beyond superficial syntactic validation. Recent applications in manufacturing show promise for validating technical documentation and ensuring consistency across engineering specifications [5].

3.4. Graph-Based Methods for Relational Data Quality

Graph-based approaches represent data as interconnected networks, providing powerful frameworks for assessing and improving quality in highly relational datasets. In sensor networks, graph-based methods offer effective means to model spatial relationships between sensors, enabling more robust anomaly detection by considering neighborhood contexts and relational constraints, with particular utility in environmental monitoring applications where spatial correlations provide valuable context for data validation [6].

4. Domain-Specific Applications and Case Studies

4.1. Healthcare: Ensuring Patient Data Integrity

The healthcare industry faces unique data quality challenges due to the critical nature of medical information and the complex ecosystem of systems that generate and consume it. Recent research in the Journal of Innovation & Knowledge shows that healthcare organizations implementing AI-powered data quality management report a significant improvement in diagnostic accuracy and treatment planning [7]. These solutions address issues ranging from simple demographic inconsistencies to critical mismatches in medication dosages, allergies, and diagnostic codes.

Electronic Health Record (EHR) systems generate massive volumes of structured and unstructured data that present significant quality challenges. Natural Language Processing techniques have proven particularly valuable for extracting and validating information from clinical notes, radiology reports, and pathology documents. The integration of AI in healthcare data management has been found to enhance clinical decision-making processes and improve overall patient care quality through higher data reliability [7].

4.2. Finance: Fraud Detection and Regulatory Compliance

The financial services sector relies heavily on data quality for risk management, regulatory compliance, and fraud prevention. According to recent studies in data analytics applications, financial institutions implementing AI-powered data quality frameworks experience significant improvements in risk prediction accuracy and fraud detection capabilities [8]. These improvements stem from the ability of machine learning algorithms to detect subtle inconsistencies and potential compliance issues before they impact official reports.

Fraud detection represents a particularly compelling application of AI-based data quality techniques in finance. Research indicates that advanced machine learning models have transformed this landscape, with predictive analytics significantly reducing false positives while simultaneously increasing fraud detection rates. The application of big data analytics in financial decision-making has been shown to substantially improve regulatory compliance while reducing operational costs [8].

4.3. Marketing: Customer Data Quality Management

Marketing effectiveness depends fundamentally on the quality of customer data, with impacts rippling through segmentation, personalization, and performance measurement. Innovation & Knowledge research demonstrates that organizations implementing AI-powered customer data quality frameworks achieve measurable improvements in customer engagement metrics and conversion rates [7]. These results stem from more accurate customer profiles, improved segmentation, and enhanced ability to track customer interactions across touchpoints.

Customer identity resolution presents a particular challenge in modern omnichannel marketing environments, where consumers interact with brands through multiple devices and channels. Studies show that organizations implementing advanced data quality management systems experience enhanced customer relationship management capabilities and more accurate customer journey mapping [7].

4.4. Manufacturing: Sensor Data Validation in IoT Environments

The industrial Internet of Things (IoT) has transformed manufacturing operations through pervasive sensing and real-time monitoring, but the massive data volumes generated by these systems create significant quality challenges. Research on data analytics transformation illustrates that manufacturing organizations utilizing advanced analytics for data validation experience substantial improvements in production efficiency and equipment reliability [8]. AI-powered quality management systems have proven essential for addressing these challenges at scale.

Time-series anomaly detection using recurrent neural networks has emerged as a particularly effective approach for sensor data validation. These techniques can identify subtle deviations from expected patterns while accounting for normal operational variations and cyclical behavior. Studies show that predictive maintenance based on validated sensor data significantly reduces unplanned downtime and extends equipment lifecycle [8].

4.5. Government: Public Records and Administrative Data Quality

Government agencies face unique data quality challenges due to the scale, diversity, and longevity of their information assets. Recent research in digital transformation indicates that public sector organizations implementing AI-based data quality systems show improved service delivery metrics and operational efficiency [7]. These approaches help address challenges while enhancing policy development and program effectiveness.

Identity resolution across government systems represents a particularly complex challenge due to fragmented data architectures and varying identifier schemes. Studies demonstrate that public sector entities implementing modern data management practices experience enhanced citizen service delivery and more effective resource allocation through improved data quality [7].

5. Challenges and Limitations

5.1. Scalability Issues with Large-Scale Datasets

As organizations accumulate ever-larger volumes of data, the scalability of AI-powered quality management systems becomes increasingly critical. Research on differential privacy and data privacy violations demonstrates that computational challenges grow significantly when implementing privacy-preserving techniques on large datasets, particularly when sophisticated inference attacks need to be detected and mitigated [9]. This creates significant performance bottlenecks for organizations attempting to maintain both data quality and privacy at scale.

Distributed computing approaches offer promising solutions to these scalability challenges, though ensuring consistency and accuracy across distributed nodes introduces additional complexity. The need to balance computational efficiency with robust privacy guarantees often leads to performance trade-offs that impact real-time quality monitoring capabilities, particularly in data intensive applications where immediate quality assessment is critical [9].

5.2. Balancing Automation with Human Expertise

While AI promises increased automation of data quality processes, finding the optimal balance between algorithmic and human involvement remains a significant challenge. Research on visual analytics and interactive machine learning demonstrates that effective human-AI collaboration requires carefully designed interfaces that make model behavior interpretable to human experts [10]. Such collaborative interfaces allow domain specialists to understand, validate, and refine algorithmic quality assessments through interactive visualization and exploration.

Effective human-in-the-loop architectures require careful design to maximize complementary strengths of human judgment and algorithmic processing. Studies show that well-designed interactive systems can significantly improve model performance by incorporating domain expertise through mechanisms such as feature selection and iterative refinement of model parameters [10]. This collaborative approach proves particularly valuable for complex quality assessment tasks requiring contextual understanding that purely automated systems struggle to achieve.

5.3. Privacy and Security Concerns

The application of AI to data quality management introduces significant privacy and security considerations, particularly when working with sensitive or regulated information. Research on differential privacy highlights the fundamental tension between data utility and privacy protection, where stronger privacy guarantees often come at the cost of reduced data usefulness for quality analysis [9]. This creates challenging trade-offs for organizations that must maintain both high data quality and strict privacy standards.

Methods for detecting potential violations of differential privacy guarantees add another layer of complexity to quality management systems. These techniques help identify when seemingly anonymized data might still leak sensitive information, but implementing them effectively requires sophisticated monitoring capabilities that many organizations struggle to maintain across their data ecosystem [9].

5.4. Model Interpretability and Transparency

As quality management increasingly relies on complex AI models, ensuring interpretability and transparency becomes both technically challenging and operationally critical. Research on visual analytics in deep learning demonstrates that understanding model behavior requires specialized visualization techniques that can reveal how models process and evaluate data quality [10]. Without such interpretability, organizations struggle to trust and effectively utilize AI-based quality recommendations.

Deep learning models present particular challenges for interpretability due to their complex internal representations. Visual analytics approaches that combine intuitive interfaces with algorithmic techniques for explaining model decisions show promise for addressing these challenges, allowing data stewards to understand why specific quality issues are flagged and how they should be addressed [10].

5.5. Handling Evolving Data Distributions and Concept Drift

The dynamic nature of real-world data presents significant challenges for AI-based quality management systems, which must adapt to evolving patterns and relationships over time. Privacy-preserving techniques face additional challenges in this context, as they must maintain their guarantees even as underlying data distributions shift [9]. This requires sophisticated monitoring and adaptation approaches that can detect potential privacy violations while accommodating legitimate changes in data characteristics.

Interactive visualization techniques offer promising approaches for identifying and responding to concept drift, allowing human experts to recognize emerging patterns and guide model adaptation accordingly [10]. These collaborative approaches leverage human pattern recognition capabilities alongside algorithmic processing to maintain quality assessment accuracy even as data evolves.

Challenge Category	Primary Challenge	Mitigation Approach
Scalability	Performance bottlenecks with large datasets	Distributed computing architectures
Scalability	Real-time monitoring limitations	Computational efficiency optimizations
Human-AI Balance	Finding optimal automation level	Carefully designed collaborative interfaces
Human-AI Balance	Complex quality assessment tasks	Interactive systems with domain expertise integration
Privacy & Security	Tension between utility and privacy	Differential privacy techniques
Privacy & Security	Sensitive information leakage	Sophisticated monitoring capabilities
Model Interpretability	Complex AI model transparency	Specialized visualization techniques
Model Interpretability	Deep learning model complexity	Visual analytics with intuitive interfaces
Concept Drift	Evolving data distributions	Sophisticated monitoring and adaptation
Concept Drift	Maintaining privacy with shifting data	Interactive visualization techniques

Table 2: Challenges in AI-Powered Data Quality Management and Mitigation Approaches [9, 10]

6. Emerging Trends and Future Directions

6.1. Explainable AI for Transparent Data Quality Processes

As AI-powered data quality systems become increasingly sophisticated, the need for transparency and interpretability has emerged as a critical requirement. Research on Explainable Artificial Intelligence (XAI) indicates that understanding how AI models make decisions is essential for building trust in automated quality assessment processes [11]. This enhanced transparency enables data stewards and business users to understand the rationale behind quality assessments, driving higher adoption rates and more effective quality improvement.

Explainable AI approaches fall into various categories including visualization techniques, textual explanations, and feature importance methods. These techniques serve different stakeholder needs, from technical teams requiring detailed model mechanics to business users who need simplified justifications for quality decisions [11]. As organizations increasingly rely on complex models for data quality management, these explanation capabilities become essential for operational acceptance and regulatory compliance.

6.2. Human-in-the-Loop Approaches for Feedback Integration

The synergistic combination of human expertise and machine intelligence continues to demonstrate superior results compared to fully automated approaches, particularly for complex quality scenarios involving ambiguous or context-dependent judgments. Research on federated learning and collaborative AI systems emphasizes that human feedback significantly improves model performance in real-world applications [12]. These hybrid approaches leverage the complementary strengths of human contextual understanding and machine scalability.

Effective human-in-the-loop quality management requires carefully designed interfaces and workflows that make efficient use of limited expert attention. By strategically directing human review to cases where algorithms have low confidence, organizations can maximize the impact of domain expertise while maintaining scalability across large datasets [12].

6.3. Transfer Learning for Cross-Domain Data Quality

Transfer learning approaches, which leverage knowledge gained in one domain to improve performance in another, are demonstrating significant promise for accelerating data quality initiatives. As explored in XAI research, these techniques allow organizations to overcome limited training data in new domains by adapting pre-trained models from related contexts [11]. This capability dramatically reduces the time and resources needed to establish effective quality monitoring for new data types.

The development of domain-adaptive quality architectures represents an emerging research direction, with frameworks designed specifically to facilitate knowledge transfer across contexts showing superior performance to general-purpose techniques. These approaches are particularly valuable for organizations that manage diverse data types across multiple business units or industry contexts.

6.4. Federated Learning for Privacy-Preserving Data Quality

Federated learning enables model training across distributed data sources without centralizing sensitive information, opening new possibilities for privacy-preserving data quality management. This approach allows multiple organizations to collaboratively train shared quality models while keeping raw data local, addressing privacy concerns and regulatory constraints [12]. For data quality applications, this means organizations can benefit from collective intelligence without exposing sensitive records.

The federated approach proves particularly valuable in regulated industries and cross-organizational contexts where data sharing faces legal or competitive barriers. As described in federated learning research, these techniques enable organizations to collaboratively improve entity resolution, anomaly detection, and other quality functions while maintaining strict data locality [12].

6.5. Real-Time Data Quality Monitoring and Remediation

The shift toward streaming data architectures and real-time analytics is driving corresponding advances in continuous quality monitoring capabilities. Explainable AI systems are increasingly important in this context, as they enable immediate understanding of quality issues as they arise [11]. This rapid detection and explanation enables timely intervention before quality issues propagate through downstream systems.

Edge computing architectures increasingly incorporate quality validation capabilities, enabling preliminary assessment and remediation before data reaches core systems. This distributed approach proves particularly valuable for IoT and sensor network scenarios where data volumes are high and connectivity may be intermittent.

6.6. Quantum Computing Applications in Data Quality

While still in early exploratory stages, quantum computing approaches offer intriguing possibilities for addressing computationally intensive data quality challenges. Similar to how federated learning enables new collaborative approaches to data quality, quantum computing may eventually transform our ability to process complex entity resolution and pattern matching tasks [12]. These emerging technologies represent the frontier of data quality research, with significant potential for addressing previously intractable quality challenges at scale.

Emerging Trend	Key Benefit	Primary Application Area
Explainable AI	Enhanced transparency and trust	Quality assessment processes
Explainable AI	Regulatory compliance	Complex model justification
Human-in-the-Loop Approaches	Improved model performance	Complex quality scenarios
Human-in-the-Loop Approaches	Efficient use of expert attention	Cases with low algorithmic confidence
Transfer Learning	Reduced training data requirements	Cross-domain quality initiatives
Transfer Learning	Accelerated deployment	New data type monitoring
Federated Learning	Privacy preservation	Collaborative model training
Federated Learning	Cross-organizational collaboration	Regulated industries
Real-Time Monitoring	Immediate issue detection	Streaming data architectures
Real-Time Monitoring	Early remediation	Edge computing environments
Quantum Computing	Advanced pattern matching	Complex entity resolution
Quantum Computing	Processing intensive challenges	Large-scale data quality issues

Table 3: Emerging Trends in AI-Powered Data Quality Management [11, 12]

7. Conclusion

The integration of artificial intelligence into data quality management represents a paradigm shift that enables organizations to address increasingly complex quality challenges at unprecedented scale and sophistication. By combining machine learning, deep learning, natural language processing, and graph-based techniques, AI systems can detect subtle patterns, relationships, and anomalies that traditional approaches would miss, dramatically improving the accuracy, completeness, consistency, and timeliness of data assets. The domain-specific applications across healthcare, finance, marketing, manufacturing, and government sectors demonstrate tangible benefits including enhanced decision-making, reduced operational costs, improved regulatory compliance,

and more personalized customer experiences. While significant challenges remain regarding scalability, human-AI collaboration, privacy protection, model interpretability, and concept drift, emerging trends in explainable AI, feedback integration, transfer learning, federated approaches, real-time monitoring, and quantum computing offer promising pathways for advancing the field. As AI continues to evolve, organizations that effectively leverage these capabilities for data quality management will gain significant competitive advantages through more reliable insights, streamlined operations, and enhanced ability to extract value from their information assets.

Funding: This research received no external funding.

Conflicts of Interest: The authors declare no conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers.

References

- [1] Amina Adadi and Mohammed Berrada, "Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI)," IEEE Access, 2018. [Online]. Available: <https://ieeexplore.ieee.org/document/8466590>
- [2] Andre Ripla, "The Digitization of the World: From Edge to Core," LinkedIn, Feb. 2023. [Online]. Available: <https://www.linkedin.com/pulse/digitization-world-from-edge-core-andre-ripla-pgcert-mqpm>
- [3] Arthur Paul Christenson Jr and William Shalom Goldstein, "Impact of data analytics in transforming the decision-making process," ResearchGate, 2022. [Online]. Available: https://www.researchgate.net/publication/366808044_Impact_of_data_analytics_in_transforming_the_decision-making_process
- [4] Fred Hohman et al., "Visual Analytics in Deep Learning: An Interrogative Survey for the Next Frontiers," IEEE, 2019. [Online]. Available: <https://ieeexplore.ieee.org/document/8371286>
- [5] J. Ding, X. Zhang, M. Hay, Y. Wang, and S. Muthukrishnan, "Detecting Violations of Differential Privacy for Quantum Algorithms," arXiv:2309.04819v1 [quant-ph] 2023. [Online]. Available: https://www.pure.ed.ac.uk/ws/portalfiles/portal/468503504/DetectingViolationsofDifferentialPrivacy_accepted_version_contribution.pdf
- [6] Jiarui Xie, Lijun Sun and Yaoyao Fiona Zhao, "On the Data Quality and Imbalance in Machine Learning-based Design and Manufacturing—A Systematic Review," Engineering, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2095809924003734>
- [7] Julius Černiauskas, "Understanding The 4 V's Of Big Data," Forbes, 2022. [Online]. Available: <https://www.forbes.com/councils/forbestechcouncil/2022/08/23/understanding-the-4-vs-of-big-data/>
- [8] Maria Trigka and Elias Dritsas, "A Comprehensive Survey of Deep Learning Approaches in Image Processing," Sensors, 2025. [Online]. Available: <https://www.mdpi.com/1424-8220/25/2/531>
- [9] Myles Suer, "What is Data Quality & Why is it Important?," Alation, 2023. [Online]. Available: <https://www.alation.com/blog/what-is-data-quality-why-is-it-important/>
- [10] Omar Ali et al., "A systematic literature review of artificial intelligence in the healthcare sector: Benefits, challenges, methodologies, and functionalities," Journal of Innovation & Knowledge, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2444569X2300029X>
- [11] Qiang Yang et al., "Federated Machine Learning: Concept and Applications," ACM Transactions on Intelligent Systems and Technology, 2019. [Online]. Available: <https://dl.acm.org/doi/10.1145/3298981>
- [12] Sarah E. McCord et al., "Ten practical questions to improve data quality," Rangelands, 2022. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0190052821000699>