

Data Quality Management Guidebook

Maximizing the Value of Data and AI

(Version 1.02)

2026-05-14

AISI Japan
AI Safety
Institute

IPA Information-technology
Promotion
Agency, Japan
Digital Infrastructure Center

Table of Contents

- 1 About This Guidebook..... 5
 - 1.1 Introduction..... 5
 - 1.2 How to Use This Guidebook 5
 - 1.3 Target Audience 6
- 2 Background and Summary 7
 - 2.1 Background 7
 - 2.2 Data Quality and Trust in Society 7
 - 2.3 AI System and Data Quality 8
 - 2.4 AI Safety and Data Quality..... 9
 - 2.5 Benefits of High-Quality Data 9
 - 2.6 Risks of Low-Quality Data 10
 - 2.7 Causes of Low-Quality Data 12
 - 2.8 Objectives and Goals..... 13
 - 2.9 Principles 13
 - 2.10 Scope 14
 - 2.11 Stakeholders in Data and AI 15
 - 2.12 Methods to Ensure Quality 16
 - 2.13 Relationship between Standards and Framework 17
 - 2.14 Quality-related Guides in Japan 17
- 3 Concept of Data Quality Management Framework 19
 - 3.1 Management, Governance and Maturity..... 19
 - 3.2 Framework Overview 20
 - 3.3 Three Views of the Framework..... 21

3.4 Data Life Cycle	22
3.5 Characteristics	23
3.6 Contents and Data Management	24
3.7 Data Quality Management for AI	26
4 Perspectives	27
4.1 Process View (Operations)	27
4.1.1 General	27
4.1.2 Data Planning	29
4.1.3 Data Acquisition	34
4.1.4 Data Preparation	38
4.1.5 Data Processing	46
4.1.6 AI System	47
4.1.7 Evaluation of Output	51
4.1.8 Delivering Results	53
4.1.9 Decommissioning	55
4.1.10 Processes throughout the Life Cycle	57
4.2 Governance Cycle View	60
4.2.1 General	60
4.2.2 Data Quality Planning	62
4.2.3 Data Quality Control	63
4.2.4 Data Quality Assurance	64
4.2.5 Data Quality Improvement	65
4.2.6 Data-related Support	66
4.2.7 Resource Provision	67
4.3 Gateway View (Characteristics)	70

4.3.1 General	70
4.3.2 Inherent Data Quality Characteristics.....	71
4.3.3 Inherent and System-Dependent Data Quality Characteristics.....	72
4.3.4 System-Dependent Data Quality	75
4.3.5 Additional Data Quality Characteristics for AI/ML	76
4.3.6 Sensor Data Quality Characteristics	79
5 Closing Remarks	84
5.1 Message	84
6 Document Information.....	84
6.1 Publishing Organizations and Contributors	84
6.2 References	85
6.3 Version History.....	85
6.4 About the Japanese Reference Translation	86
7 Appendix 1	86
7.1 Machine Learning Quality Management Guideline 4th Edition.....	86
8 Disclaimer	89

1 About This Guidebook

1.1 Introduction

Garbage in, Garbage out.

This guidebook is based on a simple premise: poor data quality directly degrades AI performance and trust. AI often amplifies existing data problems. The quality of data directly affects AI outputs.

AI systems are becoming deeply integrated into a wide range of social and industrial domains, influencing our daily lives and shaping important decisions. At the foundation of their reliability and safety lies, above all, the quality of the data they rely on. No matter how advanced an AI system may be, inaccurate or unreliable data can lead to incorrect or harmful outcomes. **Therefore, ensuring data quality is the first and most essential step in securing trustworthy AI.**

However, approaches to evaluating and managing data quality differ across organizations and sectors. In addition, the nature of data used in AI systems changes throughout its life cycle, meaning that a single, uniform standard is rarely sufficient.

This guidebook aims to provide a shared framework and way of thinking for the practical implementation of data quality management. It is informed by international standards while also incorporating practical columns and AI-specific considerations, with the goal of making it applicable and useful in real-world settings.

1.2 How to Use This Guidebook

- As guidance for policy development
 - Use this guidebook as a foundational reference or checklist when developing organizational policies for ensuring data quality.
 - Examples include data governance policies, AI ethics principles, and quality management plans.
- As a practical tool for project operations
 - Use it to clarify what needs to be checked and who is responsible at each stage of the AI development process—such as data collection, annotation, and model training.
- As a resource for education and shared learning
 - Use it as training content to help all stakeholders involved in handling data—providers, holders, developers, and users—build a common understanding.

This document is not a legally binding standard, but a guidebook intended to share practical knowledge on data quality. We hope readers will refer to it flexibly—according to their roles, needs, and purposes—and adapt it to best suit their own organizations and projects.

This guidebook refers to relevant standards, including ISO/IEC standards, and domestic and international guidelines, and organizes key concepts and checkpoints for practical consideration of data quality management in AI systems. Use of this guidebook does not directly guarantee conformity with any specific standard, compliance with applicable laws or regulations, or acquisition of any certification.

1.3 Target Audience

Data quality management cannot be completed by a single role or department; it requires collaboration across diverse functions both inside and outside the organization.

Role	Expected Uses (Examples)
Executives / Policy decision-makers	Use as guidance for establishing policies and organizational structures to ensure data quality across the entire organization and society.
Data managers / Data governance officers	Use as standards and process design references for maintaining and evaluating data quality.
Data and AI engineers	Use as practical guidance for data quality management and validation throughout the AI development life cycle.
Business planners / Service designers	Refer to it when defining and aligning necessary data quality requirements in service design and risk management.
Auditors / Researchers / Educators	Use as foundational material for evaluating, educating, and researching data quality management.

2 Background and Summary

2.1 Background

AI is rapidly transforming society by enabling more useful and efficient services. However, many people remain concerned about the use of AI, and ensuring the accuracy of AI outputs is essential for safe adoption. To achieve this, it is necessary to improve the quality of the data used for building AI and operating AI.

- AI society
 - AI systems rely on large volumes of data for training and for producing outputs such as predictions, recommendations, and explanations. The quality of data directly impacts AI system performance and effectiveness.
- Data-driven society
 - High data quality is essential in a data-driven society because it ensures accurate and reliable information. This supports informed decision-making, reduces errors, and minimizes risks. Quality data also improves customer satisfaction and helps organizations comply with regulations, protecting them from legal and financial issues.

Ensuring data quality is essential to ensuring the reliability of AI services.

2.2 Data Quality and Trust in Society

Ensuring trust is crucial for the use of AI and data. Data is the foundation of AI. If the data is unreliable, the reliability of the entire process will be compromised.

The following two lists show where data quality sits within broader discussions of trust in society and trustworthy AI.

- Elements of trust in a data-driven society
 - Service quality
 - Ethics
 - Transparency
 - Accountability
 - Privacy protection
 - Security
 - **Data quality**
 - Partnership and collaboration
- Elements of trustworthy AI

- **Accuracy and reliability**
- Ethics
- Transparency
- Accountability
- User-centric design
- Safety
- Continuous improvement

Data quality is the foundation of AI excellence, enabling trustworthy AI, which drives user adoption and engagement.

2.3 AI System and Data Quality

In AI systems, data plays an important role not only at the time of use, but also in the preparatory stage.

The figure below highlights areas related to data quality based on the OECD’s definition of an AI system. Data is an integral part of the AI system, linking the external environment and the AI model. It can be broadly categorized into input and output data used during model development and those used during operation.

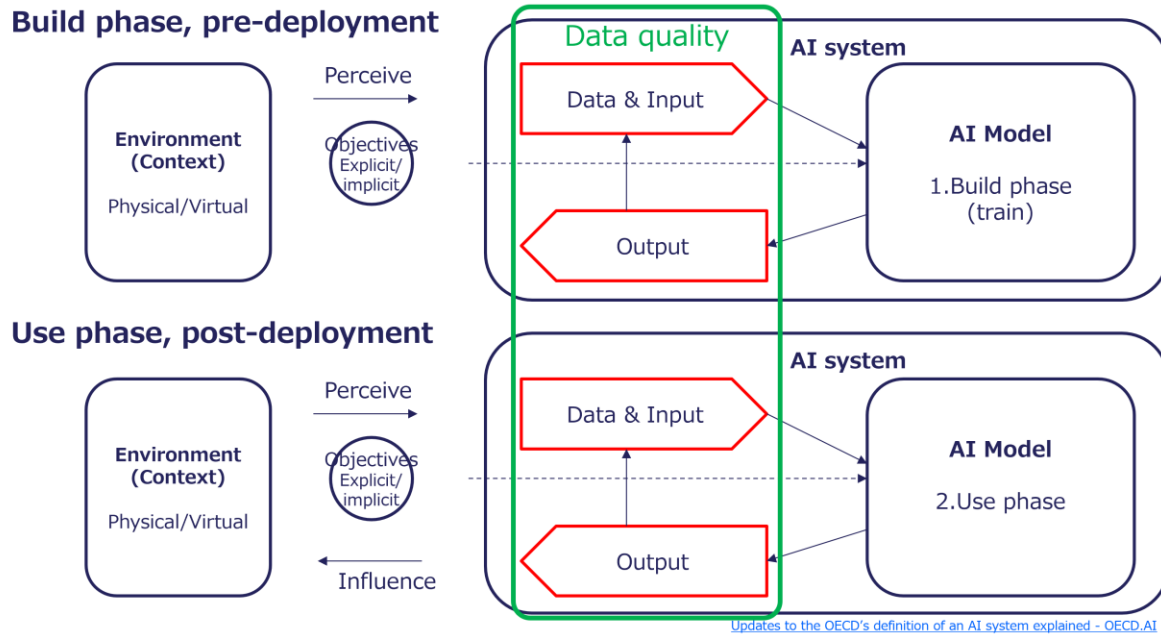


Figure 1. Relationship Between AI Systems and Data Quality

Source: OECD AI Policy Observatory (OECD.AI), 2023, “Updates to the OECD’s definition of an AI system explained”

2.4 AI Safety and Data Quality

Low-quality input data affects many critical aspects, including the performance of AI models and the credibility, consistency, and accuracy of outputs. As a result, public trust may be undermined. Data quality management is essential for ensuring the appropriate and safe use of AI. The following list organizes key evaluation perspectives related to AI safety, with data quality positioned as one of them. In addition, data quality is closely related to multiple characteristics associated with AI safety, including fairness, robustness and verifiability.

- Evaluation perspectives on AI safety
 - Control of toxic output
 - Prevention of misinformation, disinformation and manipulation
 - Fairness and inclusion
 - Addressing high-risk use and unintended use
 - Privacy protection
 - Ensuring security
 - Explainability
 - Robustness
 - **Data quality**
 - Verifiability

Source: Japan AI Safety Institute (J-AISI), 2025, [“Guide to Evaluation Perspectives on AI Safety”](#)

2.5 Benefits of High-Quality Data

Using high-quality data offers a range of benefits, including the following:

1. Enhanced accuracy
 - Higher data quality results in more accurate and consistent AI predictions and decisions, reducing the risk of flawed outcomes.
2. Better insights
 - Quality data enables the identification of actionable insights and reliable trends, empowering better decision-making processes.
3. Increased efficiency
 - Clean, high-quality data reduces time spent on preprocessing tasks such as cleaning, correction, and normalization, allowing teams to focus on higher-value activities.
4. Improved user experience

- Accurate and relevant data enhances the overall experience for end users by delivering precise, personalized, and timely results.
- 5. Reduced errors
 - Minimizing inconsistencies and errors in data reduces inaccuracies in AI outputs, improving reliability and performance.
- 6. Cost savings
 - Investing in high-quality data decreases costs related to error correction, reprocessing, and mitigating downstream issues caused by low-quality data.
- 7. Compliance and security
 - Maintaining high-quality data ensures adherence to data governance, privacy regulations, and security standards, safeguarding organizational integrity.
- 8. Enhanced trust
 - Consistently high-quality data builds trust in AI systems among stakeholders, fostering confidence in automated processes and decisions.
- 9. Scalability
 - High-quality data serves as a robust foundation for scaling AI systems, ensuring consistent performance as the system expands.
- 10. Facilitates model training and updates
 - Clean and structured data simplifies the process of training and fine-tuning AI models, accelerating development cycles and improving model adaptability.
- 11. Competitive advantage
 - Organizations leveraging high-quality data gain a significant competitive edge by delivering superior AI-driven products and services.
- 12. Ecosystem integration
 - High-quality data enables seamless integration with other systems and platforms, ensuring interoperability and streamlined workflows.

2.6 Risks of Low-Quality Data

Using low-quality data can create a range of risks, including the following:

1. Decision-making errors
 - Inaccurate or incomplete data can lead to flawed analyses, resulting in incorrect conclusions, strategic missteps, and poor decision-making at all levels of the organization.
2. Operational inefficiency

- Time and resources must be allocated to cleaning and correcting low-quality data, delaying AI model deployment and reducing operational effectiveness.
- 3. Decreased customer satisfaction
 - Inaccurate or incomplete customer data can lead to poor personalization, unmet expectations, and diminished trust in services or products.
- 4. Increased costs
 - Low-quality data increases costs through error correction, reprocessing, re-training of AI models, and potential financial losses due to incorrect predictions or recommendations.
- 5. Legal and regulatory risks
 - Non-compliance with data protection laws and regulations (e.g., General Data Protection Regulation (GDPR), California Consumer Privacy Act (CCPA)) can lead to significant fines and lawsuits.
- 6. Reputation damage
 - Errors caused by low-quality data can damage trust and credibility with customers, partners, and stakeholders, potentially leading to long-term brand erosion.
- 7. Competitive disadvantage
 - Competitors with better-quality data can outperform in key areas such as customer insights, operational efficiency, and market responsiveness, leaving organizations lagging behind.
- 8. Lost opportunities
 - Missed insights and trends due to low-quality data can result in failure to seize new market opportunities or innovate effectively.
- 9. Model performance degradation
 - AI models trained on low-quality data may exhibit bias, unreliability, or harmful behavior, leading to ethical concerns and reduced effectiveness in real-world applications.
- 10. Security risks
 - Low-quality data may inadvertently include vulnerabilities or errors that malicious actors can exploit, leading to potential security breaches or misuse of sensitive information.
- 11. Stakeholder distrust
 - Internal teams and external stakeholders may lose confidence in AI systems if the data quality consistently undermines reliability and results.

2.7 Causes of Low-Quality Data

The following are common causes of low-quality data.

- **Human error:** Operational mistakes made during data entry or data processing
- **Lack of standardization:** Inconsistent formats and standards across data sources
- **Inadequate data governance:** Lack of policies and procedures for data management
- **Outdated systems:** Using legacy systems that are not equipped to handle modern data requirements
- **Incomplete data collection:** Missing or incomplete data due to inadequate data collection processes
- **Insufficient training:** Lack of proper training for personnel handling data
- **Inadequate data integration:** Issues arising from merging data from different sources
- **Limited quality control:** Inadequate checks and validation processes for data quality
- **Unverified sources:** Using data from unreliable or unverified sources
- **Bias in data collection:** Collecting data that reflects inherent biases
- **Technical glitches:** Errors caused by software bugs or hardware failures
- **Lack of documentation (e.g., data dictionaries, data lineage diagrams):** Inadequate documentation leading to misunderstandings or misinterpretations of data
- **Misunderstandings:** Data entry errors due to misunderstanding the requirements or instructions
- **Malicious actions:** Deliberate tampering or poisoning of data by someone with harmful intent
- **Lack of incentives:** Limited motivation for end users or data contributors to provide high-quality data.

2.8 Objectives and Goals

The objectives and goals of data quality management in this guidebook are as follows:

- Objective
 - Enhancing data quality for AI ensures that decisions, predictions, and recommendations are based on accurate, consistent, and reliable information. **High-quality data empowers organizations to leverage AI effectively, reduce risks, and unlock its full potential.** At the societal level, **improving data quality helps foster trust in AI systems, promote ethical data use, and support fairness**, thereby contributing to innovative applications in healthcare, education, and public services, as well as **to economic growth and a better quality of life for all.**
- Goal
 - **The goal of data quality management is to establish standardized practices** for data collection, cleaning, validation, and management, thereby **creating an environment in which accurate, secure, and accessible data enables AI systems to achieve their best performance.** At the societal level, the aim is to broaden the benefits of data-driven innovation across sectors, promote transparency, enhance decision-making processes, and accelerate progress toward a sustainable and inclusive digital society where technology serves humanity responsibly and effectively.

2.9 Principles

This guidebook considers the following as the principles of Data Quality.

- **Know your users and their needs.**
- **Ensure that it is fit-for-purpose.**
- **Consider trade-offs between quality and cost.**
- **Accept that errors will occur, don't pursue perfection.**
- **Review the design and consider the life cycle.**
- **Use mature services where suitable capabilities already exist.**
- **Get feedback from all stakeholders.**
- **Visualize for easy confirmation and traceability.**

The required level of data quality depends on whether the data is fit-for-purpose. In the context of AI safety, this way of adjusting requirements according to the intended use, risks, and potential impact of the AI system is often described as a risk-based

approach. In high-risk use cases, provenance, representativeness, bias, personal data, output evaluation, and auditability should be checked with particular care.

Quality and cost are also important. When planning for quality, costs must be considered.

- Higher quality usually costs more. (a)
 - If you raise the quality above the target, it will cost more.
- It is important to incorporate quality in early processes such as design. (b)
 - If the design is appropriate, operating costs and total costs will decrease, and you will be able to respond quickly to changes.
- Consider not only initial costs but also operational costs. (c)
 - If you optimize your data when you implement the system, your operating costs can be reduced.

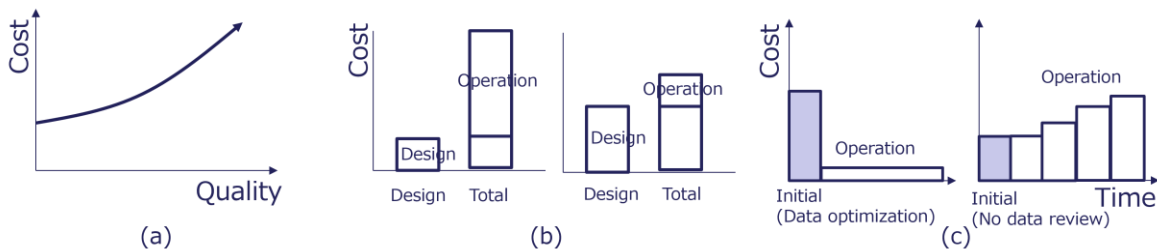


Figure 2. Quality and Cost

2.10 Scope

AI and data play important roles in society, but their scope is broad. They therefore need to be considered within a wider social context.

The figure below shows how data and AI systems are positioned within a broader hierarchy extending from systems to human interaction and social impact. This guidebook focuses primarily on the domain of AI systems, including data and data systems, within this wider context.

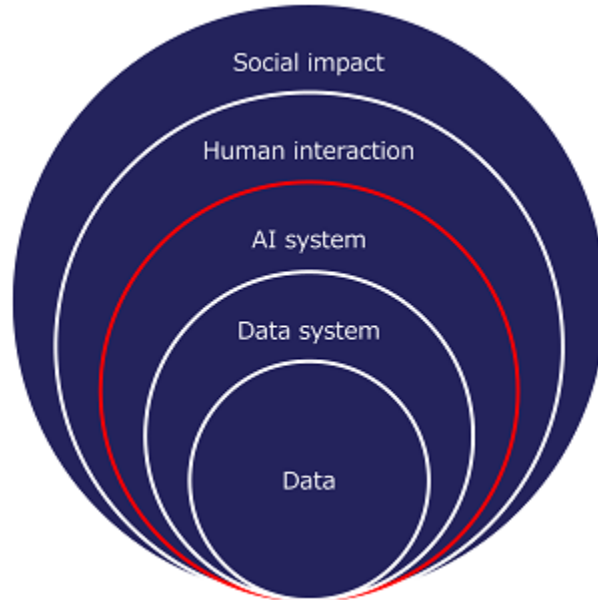


Figure 3. Scope and Positioning of AI and Data

AI systems handle a wide variety of data, including structured, semi-structured, and unstructured data such as numerical/tabular data, text, images, audio, and video. This guidebook presents general approaches that do not depend on any specific data type. AI systems include both conventional task-specific AI systems and generative AI systems, but this guidebook focuses on concepts that apply to both without drawing a strict distinction between them. That said, it also includes some discussion specific to generative AI systems, such as Retrieval-Augmented Generation (RAG).

2.11 Stakeholders in Data and AI

Stakeholders are diverse. They are characterized by being both users of AI and suppliers of data.

The figure below illustrates the actors involved in the data and AI domains and the relationships among them. A wide range of stakeholders are interconnected, with AI and data at the center.

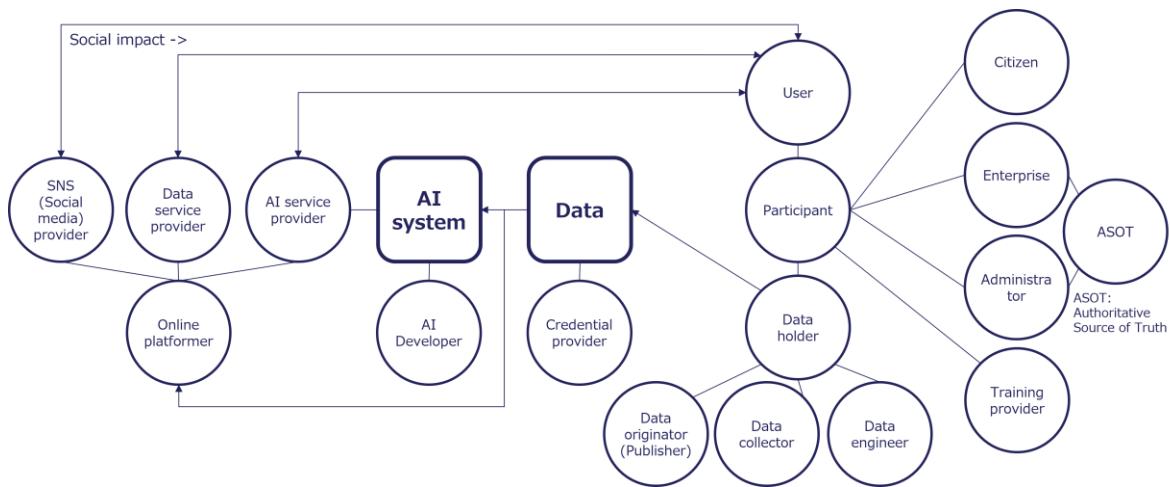


Figure 4. Stakeholder Relationship Map for Data and AI

2.12 Methods to Ensure Quality

The following methods are the main approaches to ensure data quality.

- **Make high-quality data**
 - To ensure advanced data design, utilize appropriate reference models and data modelers. Avoid manual intervention during data creation and perform validation at each stage.
- **Prevent low-quality data**
 - Use validators, detectors, converters, and cleansing tools to validate input and output data. Also perform human checks.
- **Build trust**
 - Verify the provenance of the data. Utilize digital content transparency technologies. Embed watermarks.
- **Use it right**
 - Ensure appropriate use of data in AI training. Take UI/UX into consideration. Additionally, manage data using the Single Source of Truth (SSOT) and Authoritative Source of Truth (ASoT) concepts. Also, label the data appropriately.
- **Establish governance**
 - Implement appropriate governance, such as by utilizing the DataOps approach.

2.13 Relationship between Standards and Framework

The framework provided in this guidebook organizes and guides the practical application of many data quality standards, including ISO and IEC standards, in a way that makes them easier to implement. It does not create new standards, but is only a practical guide, and synergies can be achieved by feeding the results of implementation back into the standardization process.

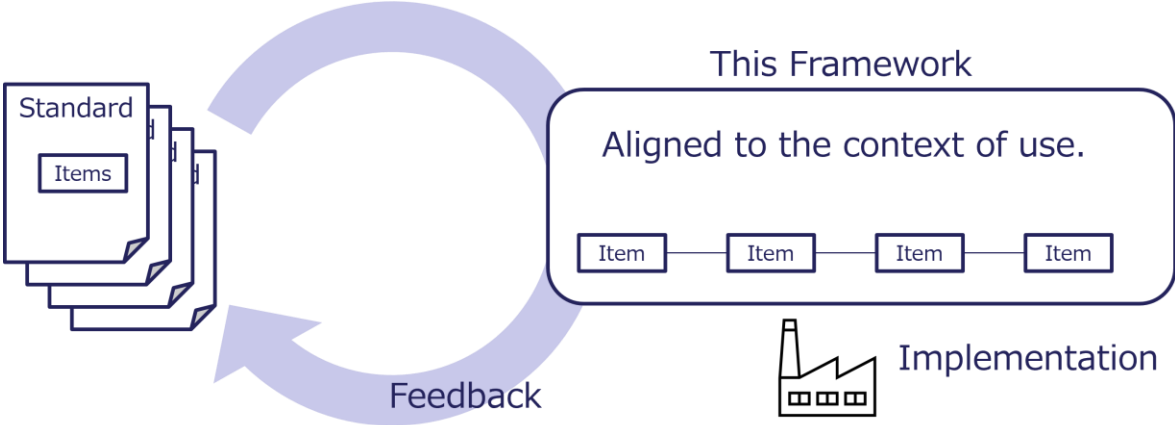


Figure 5. Framework for Applying Standards in Practice

2.14 Quality-related Guides in Japan

In Japan, several guides concerning AI quality have been established.

The figure below presents a hierarchical overview of the major AI quality-related guides in Japan. It first provides an overall perspective through the “AI Guidelines for Business”, which aims to ensure AI safety. Next, appropriate guides are selected according to the purpose of AI use and the characteristics of the AI model. These AI systems are supported by high-quality data prepared using the Data Quality Management Guidebook. In addition, frameworks on data design and interoperability underpin these efforts to enhance data quality.

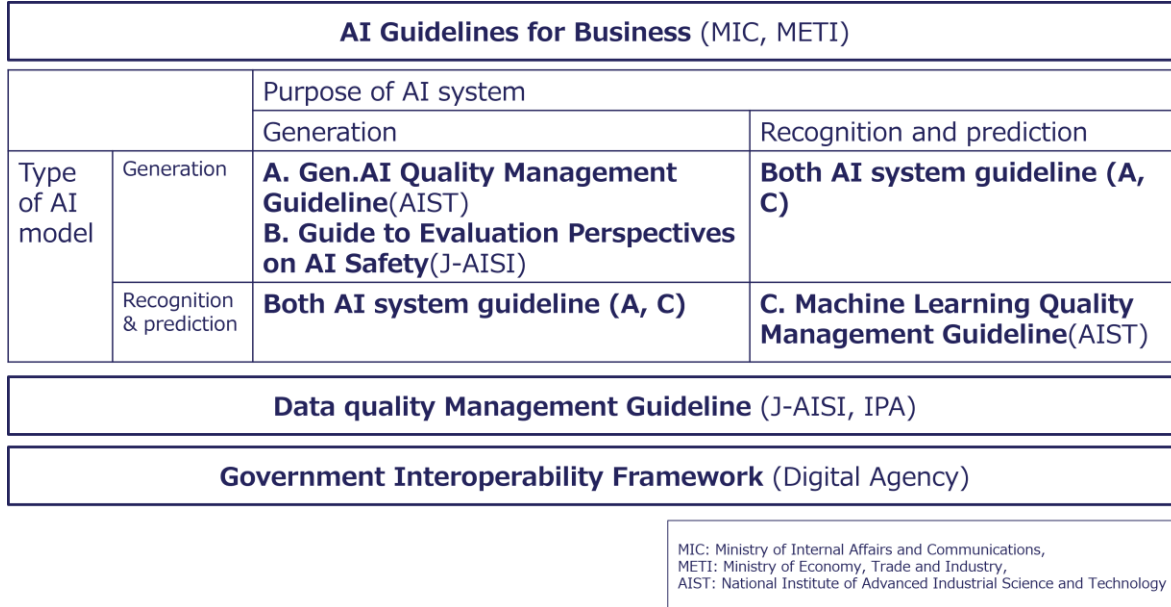


Figure 6. Structure of AI Quality Guides in Japan

The figure below illustrates the relationships shown in the previous Figure from a data perspective.

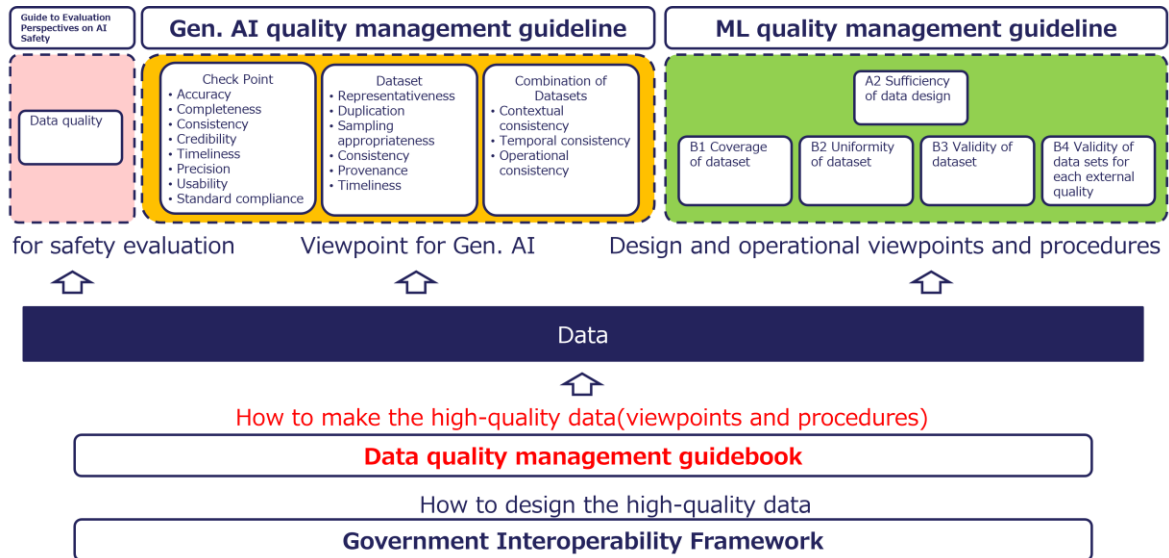


Figure 7. Relationships among AI Quality Guidelines and the Role of Data Quality Management

Japan’s high service quality is realized based on the quality of its products, people and data.

The figure below illustrates how products, people, and data interact to support high-quality services. Both service providers and users share a strong commitment to quality, and their actions are reflected in the quality of data and products, forming a continuous cycle that enhances service quality.

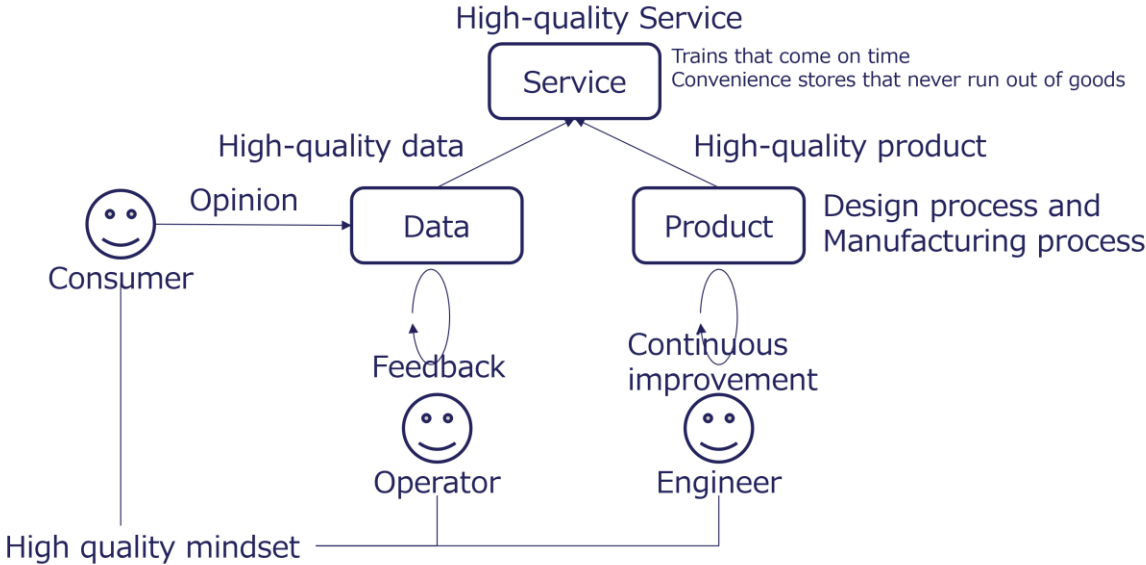


Figure 8. Structure of High Service Quality in Japan

3 Concept of Data Quality Management Framework

3.1 Management, Governance and Maturity

Continuously providing high-quality data is essential to promote the widespread adoption of AI. Effective management and governance are key to maintaining data quality and achieving maturity in data and AI utilization.

1. Data and AI Management
 - Appropriate data and AI management ensures data and AI quality.

*This guidebook mainly focuses on this area.

2. Data and AI Governance

- Appropriate data and AI governance ensures sustainable data and AI quality.
3. Data and AI Maturity
- If quality-assured data is supplied, the use of data and AI will expand.

The figure below illustrates the relationship between management, governance, and maturity. First and foremost, value is created through AI by ensuring that data is properly managed and that governance and maturity are maintained. This guidebook focuses first on management as the initial step.

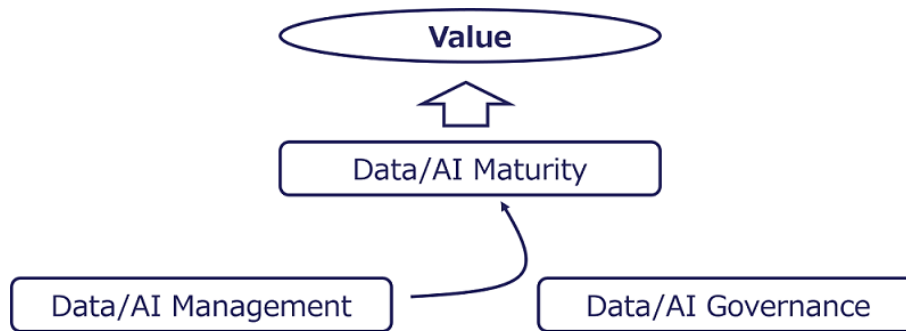


Figure 9. Relationship Between Management, Governance, and Maturity

3.2 Framework Overview

Because data is exchanged across systems, organizations, and borders, interoperability is essential.

There are numerous international standards and industry guidelines related to data quality. This guidebook integrates them into three categories: characteristics, processes, and governance.

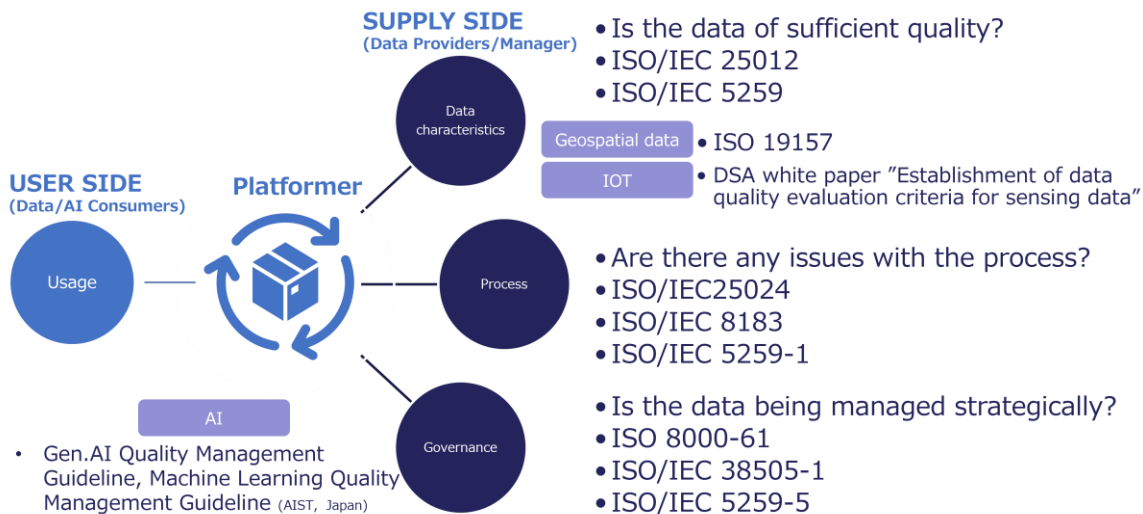


Figure 10. Organization of International Standards for Data Quality

3.3 Three Views of the Framework

This guide organizes data quality from multiple perspectives. It uses three views: the **Process View**, the **Governance Cycle View**, and the **Gateway (data quality characteristics) View**.

The **Process View** focuses on activities related to data. It covers the full set of activities such as planning, acquisition, processing, and use of data. It also includes upstream activities such as design and requirements definition.

The **Governance Cycle View** focuses on the mechanisms used to execute, sustain, and improve these activities. It addresses organizational aspects such as policies, roles, responsibilities, evaluation, and continuous improvement. Through this view, data quality is managed at the organizational level.

The **Gateway (data quality characteristics) View** focuses on the data itself. It evaluates the resulting data quality in terms of characteristics such as accuracy, consistency, and completeness. It also includes characteristics related to data description and format, such as portability and traceability.

Each of these three views has a distinct role:

- **Process View:** Activities (how quality is created)
- **Governance Cycle View:** Mechanisms (how activities are managed and sustained)
- **Gateway View:** Data State (what quality is achieved)

By combining these three views, it becomes possible to manage evaluation, execution, and continuous improvement in an integrated manner.

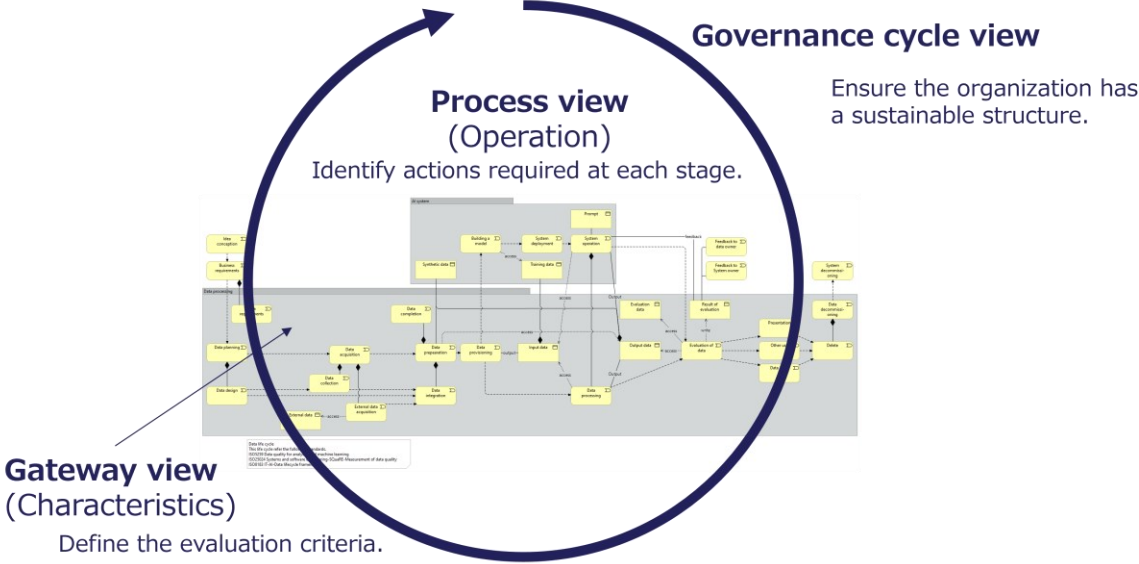


Figure 11. Three-View Framework for Data Quality Management

3.4 Data Life Cycle

In data quality management, attention tends to focus on downstream processes such as data processing and integration. These processes assume that data already exists. However, actual data quality is significantly influenced by upstream activities, including design and data acquisition.

Data should be managed not only during its utilization phase. It is also important to consider its eventual disposal. This helps prevent unnecessary data accumulation and ensures appropriate data management.

Therefore, this guide does not treat data as a one-time processing target. Instead, it treats data as part of a continuous flow from design through to disposal. This flow is referred to as the data life cycle.

Various standards define the data life cycle. In this guide, multiple standards are integrated and organized. Based on this, a detailed data life cycle is defined as illustrated in the figure below.

- Unclear data
- **Data Quality Characteristics (ISO/IEC 25012 and ISO/IEC 5259-2)**
 - Accuracy
 - Completeness
 - Consistency
 - Credibility
 - Currentness
 - Accessibility
 - Compliance
 - Confidentiality
 - Efficiency
 - Precision
 - Traceability
 - Understandability
 - Availability
 - Portability
 - Recoverability
 - Auditability
 - Balance
 - Diversity
 - Effectiveness
 - Identifiability
 - Relevance
 - Representativeness
 - Similarity
 - Timeliness

3.6 Contents and Data Management

AI handles various data, and its stakeholders are diverse. This diagram illustrates the multilayer structure of the AI and data ecosystem in relation to data quality, transparency, and risk management. This diagram is intended to provide an overview, showing how data quality, transparency, and governance relate across the ecosystem.

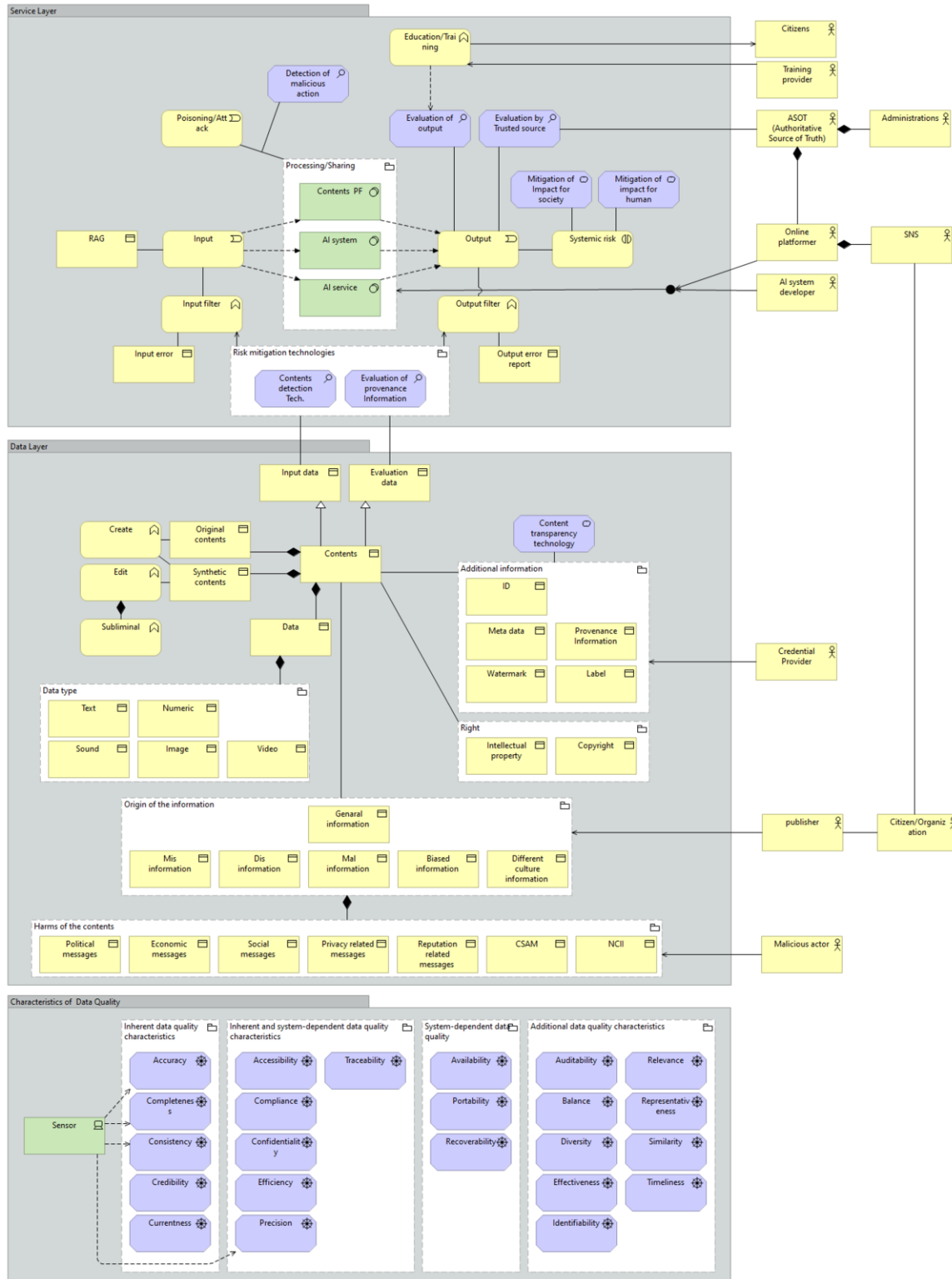


Figure 13. Integrated Structure of Contents to Be Managed for AI System

3.7 Data Quality Management for AI

When managing data quality, accuracy, completeness, and timeliness are important. However, when data is used for AI, additional considerations are required, including the following:

1. Annotation (labeling) quality management
 - For machine learning, manual or automated labeling of data is essential. If there are many incorrect or inconsistent labels, it can lead to reduced model accuracy and skewed training results.
2. Checking for bias and ensuring fairness
 - It is essential to verify that the data does not disproportionately represent certain attributes (e.g., gender, race, region) and to prevent such biases from being amplified during the training process.
3. Privacy and security measures
 - When handling data that includes personal or confidential information, proper anonymization and masking must be performed, and secure storage and processing must be ensured. Additionally, privacy protection technologies (e.g., differential privacy) should be considered to prevent personal information from being inferred using trained models.
4. Version control and drift management
 - AI models are trained based on the data distribution at the time of training. Over time, the data distribution—or the underlying relationships between input and output—may change (“data drift” or “concept drift”), which can cause a sudden drop in model performance. It is important to implement data version control, monitor distributions, retrain or refine models as needed.
5. Handling synthetic data and generated content
 - In cases where real-world data is insufficient, synthetic data may be used. However, clarity regarding its generation and quality is critical. Improper use of synthetic data can lead to flawed training and incorrect model outcomes.
6. Ensuring explainability and transparency
 - AI models should not operate as “black boxes.” Explainable AI (XAI) methods help clarify the decision-making process behind critical outcomes, thereby enhancing trust and accountability.
7. Establishing operational and monitoring frameworks
 - Continuous monitoring and maintenance are essential not only during data and model development but also throughout deployment. It is

important to have systems in place that can respond quickly if unexpected bias or performance decline arises.

4 Perspectives

4.1 Process View (Operations)

4.1.1 General

The Process View considers data quality from the perspective of whether the necessary activities are appropriately carried out along the data life cycle. Therefore, this view evaluates the processes that operate data, rather than the data itself.

These processes include both those that correspond to stages in the data life cycle and those that are performed across multiple stages.

In this guide, the life cycle is divided into the following eight major stages. The processes are then organized in correspondence with each stage.

1. Data Planning
2. Data Acquisition
3. Data Preparation
4. Data Processing
5. AI System
6. Evaluation of Output
7. Delivering Results
8. Decommissioning

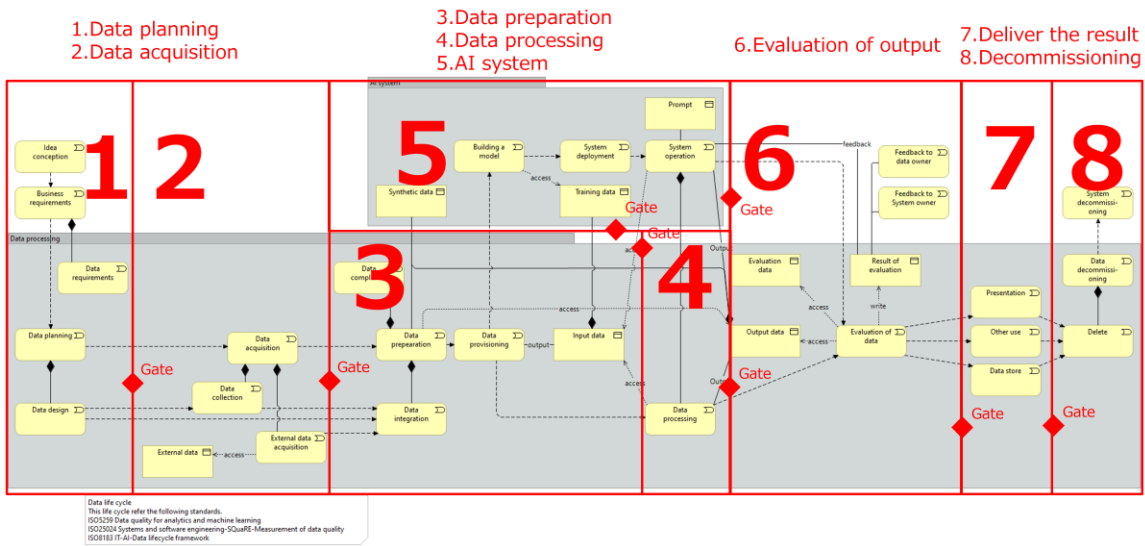


Figure 14. Data Life Cycle Processes

4.1.1.1 Role of Stakeholders

The stakeholder map below shows which actors are involved across the process. This helps clarify role allocation and collaboration across organizations.

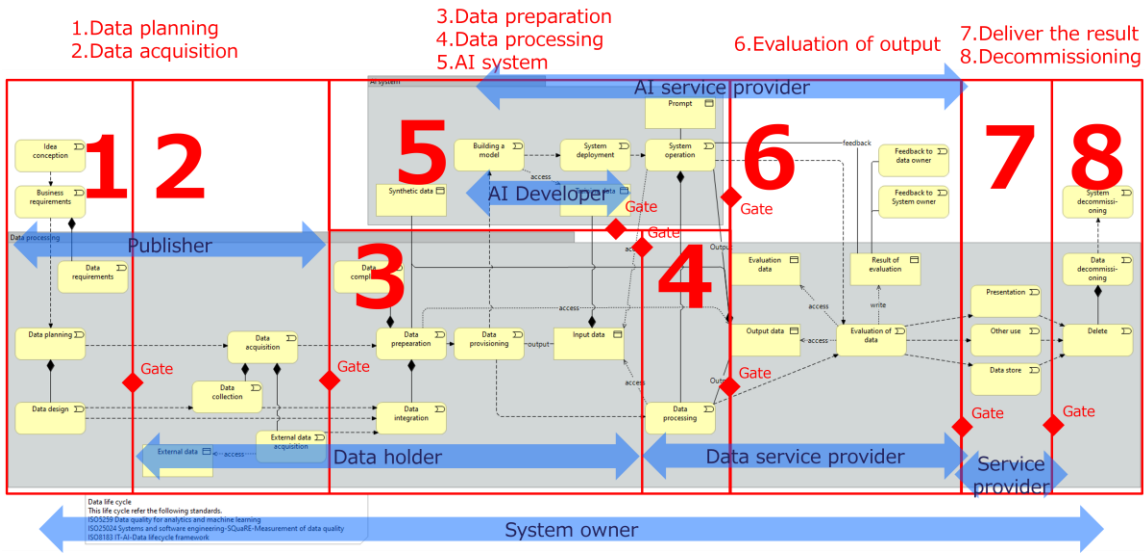


Figure 15. Roles in the Data Life Cycle

4.1.1.2 Structure of This Section

The following sections describe eight processes corresponding to the stages, as well as processes that are carried out across the data life cycle. Each is organized using the structure below.

- **General**
 - Describes the purpose and role of the process. It also outlines the main issues, requirements, and expected value. The description of a process may include an overview of the stage to which the process corresponds.
- **Processes**
 - Presents the individual processes that constitute the process. For each process, example procedures and checkpoints for evaluating implementation are provided. Checkpoints may refer to the state of data for simplicity. However, in the Process View, the focus is on whether these states are appropriately verified and managed. It is assumed that, in practice, the procedures and checkpoints will be adapted as appropriate depending on the context and environment.

4.1.2 Data Planning

4.1.2.1 General

- **Description**
 - Data planning is a critical process in the data life cycle.
 - This process covers the definition of the intentions and motivations behind the need for the data.
 - It includes planning for interoperability across the entire service and scalability for future growth.
 - It specifies the data life cycle, including acquisition methods, evaluation methods, and disposal processes.
- **Issue**
 - Difficulty in managing and organizing existing data effectively
 - Insufficient revision and updating of rules and processes to align with current needs
 - Difficulty in leveraging new technologies to optimize data management practices
- **Requirement**
 - Processes for defining and managing high-quality data, including standardized and well-structured formats
 - Mechanisms for ensuring interoperability between systems and services

- Processes for selecting and maintaining a reliable and authoritative “source of truth” for data consistency
- Value
 - Enhanced efficiency in downstream processes through better-prepared data
 - Greater sustainability and scalability of the overall service architecture

Quality is built through three steps: concept development, requirements definition, and design.

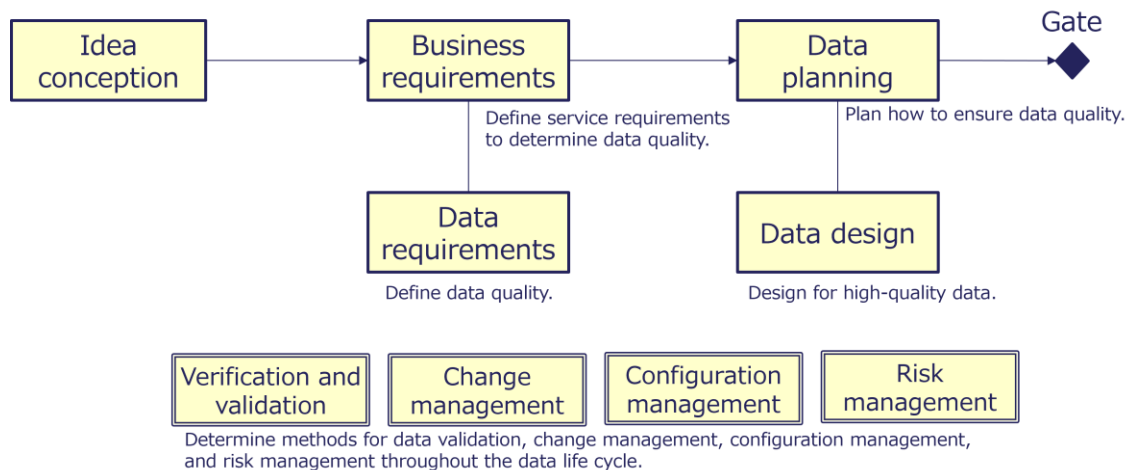


Figure 16. Data Planning Process

4.1.2.2 Idea Conception

- Procedure
 1. Gather user needs.
 2. Create a concept.
 3. Define the business and data policy.
 - Align with the organizational policies.
 - Check the future trend.
 4. List the stakeholders.
 5. Check the feasibility.
- Checkpoint
 - Do decision-makers understand the overall system concept?
 - Do decision-makers understand the benefits of high-quality data?
 - Do decision-makers understand the risks posed by low data quality?

- Do decision-makers understand the costs of requiring higher data quality than necessary and the importance of balancing effectiveness and cost?
- Do you have organizational capacity and skilled personnel?

4.1.2.3 Business Requirements

- Procedure
 1. Define the business objective and scope.
 2. Define the value, requirement, constraint, and risk.
- Checkpoint
 - Is service quality clearly defined? (speed, cost, etc.)

4.1.2.4 Data Requirements

- Procedure
 1. Define the relevant data.
 - Input, Output, Reference
 - Constraint, Interface, Statistical data
 2. Specify the data specifications.
 - Description, Goal, Requirement, Matters to consider
 - Refer to ISO/IEC 5259-3 for data specification.
 3. Define the required quality level.
- Checkpoint
 - Are the data required by the business listed?
 - Have quality control items been defined for each data?
 - Are the data quality requirement levels for each data defined?

4.1.2.5 Data Planning

- Procedure
 1. Check the existing data.
 2. Gather the data needs.
 3. Define the master data.
 4. Define the data architecture.
 5. Define the design policy and methodology.
 - Structure, Location, Modeling, Documentation
 6. Define and find the data source.
 7. Check whether any data-related legislation applies.
- Checkpoint

- Do decision-makers agree on the conversion of existing data to the new model?
- Are the data needs of stakeholders for their operations and decision-making understood?
- Are the data architecture and design policies documented?
- Are the required data defined and available?
- Have legal constraints on the use of the data been checked?
- Does it contain any personal, sensitive, or privacy-related data?
- Does it contain any intellectual property data?

4.1.2.6 Data Design

- Procedure
 1. Check the reference models and taxonomy.
 2. Design the data models and the taxonomy.
 3. Design the metadata and labels.
 4. Design the rules.
 5. Check compliance with applicable legislation.
 - Note: In Japan, GIF (Government Interoperability Framework) provides reference models.
- Checkpoint
 - Do you refer to a data reference model or standardized taxonomy?
 - Do you use modeling tools?
 - Is the metadata designed based on a Data Catalog Vocabulary (DCAT)?
 - Do you refer to general rules for utilization and access?
 - Are the data rules and models designed to comply with applicable legislation?

Column: Reference Model

A reference model is a conceptual framework that provides standardized structures, processes, and data frameworks for a specific industry or domain. It serves as a “template” for designing and implementing systems or processes, offering unified terminology, definitions, and methods. This enhances interoperability between organizations and systems, improving efficiency in development and operations. Due to its high level of abstraction, a reference model can be applied to a wide range of scenarios and customized to meet specific requirements.

- How a Reference Model Improves Data Quality

- A reference model establishes a foundation for data consistency and standardization.
- By adopting common data definitions and naming conventions, it prevents inconsistencies across different systems.
- Additionally, using a reference model eliminates data duplication and redundancy, enabling the creation of an efficient and streamlined data structure.
- By embedding governance rules and constraints aligned with business processes, it enhances the accuracy and integrity of the data.

Reference models play a critical role in ensuring data completeness during the design phase and reducing errors and costs of correction during the operational phase.

Column: Data Dictionary

A Data Dictionary is a repository that systematically organizes and manages the names, definitions, formats, and relationships of data elements used in information systems and databases. It defines attribute names, data types, lengths, constraints, and business rules for each data element, providing consistent standards across the organization.

This shared understanding of data meaning and structure among stakeholders—such as developers and operators—helps streamline system development and operations, improving efficiency and quality. Moreover, during data mapping and system integration, it enhances the accuracy of data migration and reporting between different systems, making it easier to ensure consistency.

In terms of data governance, a Data Dictionary serves as the foundation for security requirements, access rights, and regulatory compliance.

By keeping it up to date, organizations can flexibly accommodate new data requirements and swiftly adapt to changing business needs, ensuring a robust and responsive system environment.

Column: Data Modeling

Data modeling is the process of visually and logically organizing the data used in information systems and databases, representing the structure and relationships of the data in the form of a model. It involves defining entities (data objects), attributes (data elements), and relationships (connections between data). Data modeling

typically includes three levels: conceptual data models (representing data from a business perspective), logical data models (defining detailed structures from a technical perspective), and physical data models (defining the database structure concretely).

Data modeling contributes to improved data quality as follows.

1. Ensures consistency and standardization
2. Eliminates duplication and redundancy
3. Ensures data integrity and consistency
4. Enhances data reusability
5. Improves data-driven decision-making

Data modeling is not just a technical process; it is a critical activity that underpins data quality management and data governance.

Column: DCAT

DCAT (Data Catalog Vocabulary) is a W3C standard vocabulary designed to facilitate interoperability between data catalogs published on the internet. It enables publishers to describe datasets and data catalogs in a standardized way, making them easier to find, share, and reuse across platforms. DCAT supports metadata such as titles, descriptions, URLs, and licensing. By promoting consistent metadata, DCAT improves the discoverability of datasets and supports data integration and open data initiatives. The current version, DCAT v3, extends support for data services, versions, and provenance, aligning with FAIR data principles.

- FAIR: Findable, Accessible, Interoperable, Re-usable

Source: World Wide Web Consortium (W3C), 2024, [“Data Catalog Vocabulary \(DCAT\) - Version 3”](#)

4.1.3 Data Acquisition

4.1.3.1 General

- Description
 - Data acquisition involves generating, collecting, and acquiring data from internal and external sources.
 - This process includes measures for preventing errors during data acquisition, whether from sensors or manual data input.

- It includes verification activities to exclude low-quality data when obtaining data from external sources.
- When supplementing real-world datasets with synthetic data, this process includes documentation and controls to prevent the inclusion of inappropriate synthetic content.
- It includes adding metadata, such as data profiles and provenance information, to support traceability and data quality management.
- Issue
 - Risk of including low-quality or inappropriate data during data collection
- Requirement
 - Processes for collecting high-quality data efficiently and minimizing inaccuracies
 - Measures for maintaining quality at the point of acquisition and reducing the need for extensive pre-processing steps, such as data cleansing
- Value
 - Reduced time and costs associated with data preparation and cleaning
 - Increased trust in the data, leading to better decision-making and enhanced credibility for AI applications

Generate and collect the data required to achieve the objective, and if the data is not available within the organization, acquire it from external sources.

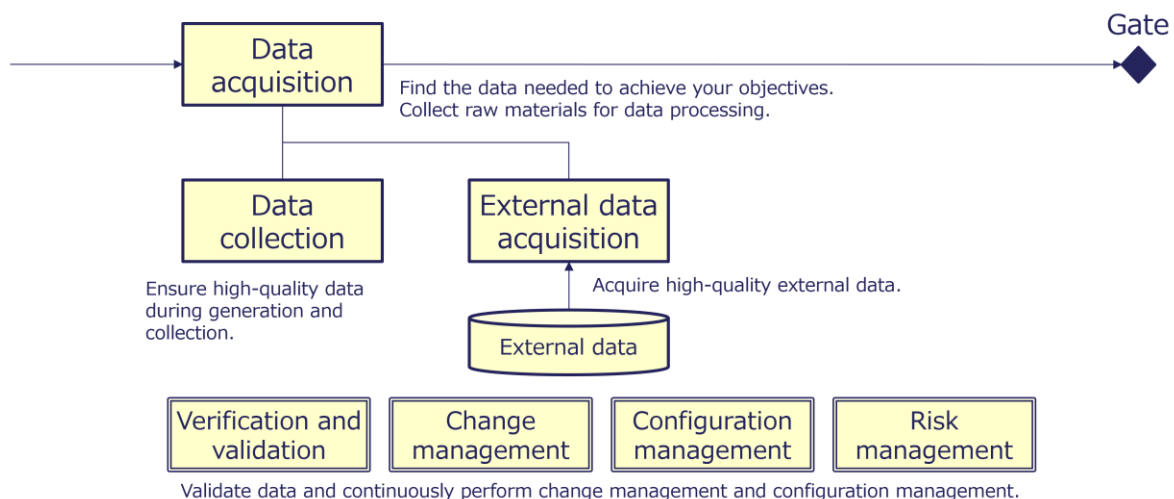


Figure 17. Data Acquisition Process

4.1.3.2 Data Acquisition

- Procedure
 1. Identify and locate necessary data.
 2. Check the provenance information.
 3. Check the condition for use.
- Checkpoint
 - Is data obtained from reliable sources?
 - Are there any problems with the data's provenance information?
 - Are there any restrictions on how the data can be used?

4.1.3.3 Data Collection

- Procedure
 1. Check the device (Sensor).
 2. Collect or input the data.
 - Prevent errors by using the web forms and APIs.
 3. Verify the data.
 - Remove out-of-range and inappropriate data.
 - Correct inconsistent data.
 4. Anonymize and mask.
 5. Create metadata.
 - General information
 - Quality information
 - Method of measuring or gathering data
- Checkpoint
 - Are you taking steps to prevent inappropriate data?
 - Do you ensure that out-of-range data is not entered?
 - Do you check that the data is consistent?
 - Do you refer to DCAT for metadata?
 - Do you describe the method of measurement and collection of data?

4.1.3.4 External Data Acquisition

- Procedure
 1. Verify the metadata and provenance information.
 2. Verify the data.
 - Checking quality characteristics
 - Finding inappropriate content (including synthetic data)

3. Add metadata and provenance information.
- Checkpoint
 - Is the source of the data reliable?
 - Is the data clear about provenance information?
 - Does the quality of the data meet the specifications?
 - Does the data meet the requirements for “Portable Data”, allowing for effective exchange and transfer with the target system without loss of content or meaning?
 - Does it contain data inappropriate to the system’s objectives?
 - Does the data contain important information other than the data items to be imported (such as Personally Identifiable Information (PII))?
 - Does it contain synthetic data that is not clearly identified as synthetic data?
-
-

Column: Improvement of Input Method

By linking data via an interface between systems, we can prevent input errors by users. The following effects are expected as a result.

- Accurate and quick
- No delivery errors

When a person is entering data, using an input form allows the user to correct or adjust the data as they enter it. Specifically, it facilitates the following:

- Checking the data formats
 - Checking the accuracy of the data
 - Preventing input variations
 - Using controlled vocabulary
-
-

Column: Provenance Data

Provenance data refers to information that records the history and origin of data, including its creation, transformation, movement, and usage. It tracks who created, modified, or used the data, as well as when and how these actions occurred. Provenance data serves as metadata to enhance transparency in data processing, making it particularly valuable in fields such as healthcare, finance, and scientific research where data reliability and integrity are critical. This information typically

includes the data's source, processing history, applied business rules or algorithms, and associated systems or users.

- How Provenance Data Improves Data Quality
 - Provenance data ensures data reliability by making its history transparent.
 - By clarifying the origins and life cycle of the data, it becomes easier to validate its accuracy and appropriateness, and to identify errors or tampering.
 - It also provides insights into data creation and updates, enabling the correction of inconsistencies or incomplete records.
 - Additionally, in audits and compliance processes, provenance data serves as evidence that the data has been handled appropriately, boosting confidence in its reliability.

This enhances the accuracy and effectiveness of data analysis and decision-making, while also mitigating potential risks and strengthening data governance.

4.1.4 Data Preparation

4.1.4.1 General

- Description
 - Data preparation is the final process for making data available at a quality level suitable for use.
 - This process includes cleaning prepared data by removing erroneous or out-of-range data and adjusting taxonomy conversions, semantics, and granularity.
 - It includes supplementing missing data as necessary.
 - It includes integrating cleaned data into a unified dataset.
 - It includes adding watermarks to the data, when necessary, to help support authenticity verification and discourage misuse.
 - It may include creating additional data for data augmentation using machine learning techniques.
- Issue
 - Need for refinement and standardization of data definitions and quality
 - Unclear or incomplete data profiles
 - Insufficient data quality management or missing key data elements
- Requirement
 - Processes for ensuring data quality appropriate for intended use

- Processes for creating comprehensive and accurate metadata about the data
- Management of provenance information to trace the source and transformation history of the data
- Measures for assuring data authenticity and reliability
- Value
 - More accurate and reliable processing through the preparation of high-quality input data
 - Improved data validation, confirmation, and traceability throughout the data life cycle

Various datasets are integrated to create input data, which is used for training AI models or further processing.

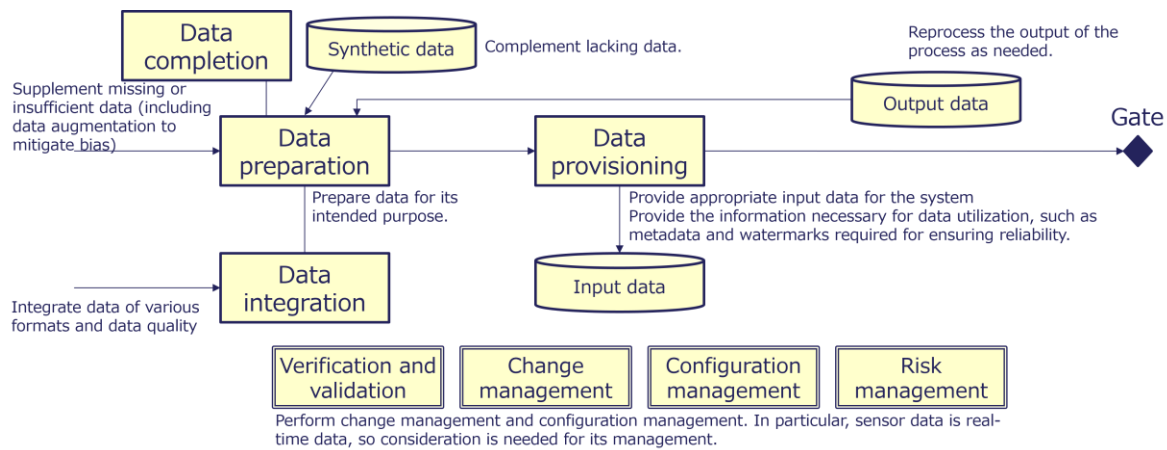


Figure 18. Data Preparation Process

4.1.4.2 Data Preparation

- Procedure
 1. Document the collected data.
 2. Define the data integration policy.
 3. Define the data supplementation policy.
 4. Clean the data.
 5. Add labels.
- Checkpoint
 - Is all necessary data documented?
 - Is a data integration policy in place?
 - Is there a defined policy for data supplementation?

- Is the dataset split logically for AI training and validation?

4.1.4.3 Data Integration

- Procedure
 1. Define the post-integration data model.
 2. Create conversion tables for taxonomies and controlled vocabularies.
 3. Verify the semantics of data items and determine the matching method.
 4. Harmonize accuracy and units.
 5. Convert file formats (e.g. XML, JSON, CSV).
 6. Document semantic differences.
 7. Create integration-related metadata.
- Checkpoint
 - Has the conversion table been created taking into account the meaning of the data items?
 - Are there priority rules for duplicate data?
 - Is blank data treated as no data?
 - Is data conversion carried out automatically by tools?
 - Is the consistency of integrated data checked?
 - Is metadata created in accordance with DCAT?

4.1.4.4 Data Completion

- Procedure
 1. Analyze the state of the data and determine whether it needs to be supplemented.
 2. Add complementary data.
- Checkpoint
 - Are you using appropriate data to supplement the data?

4.1.4.5 Data Augmentation (for AI system)

- Procedure
 1. Prepare the base data.
 2. Generate training data based on the base data.
- Checkpoint
 - Is there any bias in the base data?
 - Is there a bias caused by reliance on a small number of base data?

4.1.4.6 Synthetic Data Acquisition

- Procedure
 1. Check the metadata and provenance information.
 2. Validate the data.
- Checkpoint
 - Is it clearly indicated that the content is synthetic?
 - Does it contain inappropriate or non-consensual content?

4.1.4.7 Data Provisioning

- Procedure
 1. Register catalog information.
 2. Provide data interface.
 3. Control the version.
 4. Provide data samples.
 5. Add content protection information, such as watermarks.
 6. Manage access controls.
 7. Track the usage, if necessary.
- Checkpoint
 - Is the data provided in a machine-readable interface?
 - Do you include mechanisms such as watermarking to make it difficult for data to be misused?
 - Are access controls in place to prevent unauthorized use?

Column: IDs

Assigning unique identifiers (IDs) to data is critical for improving data quality. By tagging each record with an ID, you make it possible to integrate information from multiple sources and ensure consistent data management across systems. This prevents confusion or duplication when similar or related pieces of data need to be merged or compared. For example, if two departments track customers in separate databases, using the same customer ID allows the two datasets to be combined accurately.

In addition to aiding data integration and consistency, assigning IDs provides several other benefits:

1. Traceability and auditability: A unique ID makes it easier to trace the origin and history of a data record, which helps in audits and compliance checks.
2. Efficient updates: When you need to update or delete specific data, having an ID allows quick and precise identification of the target records.
3. Scalability: As data grows, unique IDs ensure that new records can be added without conflicts or confusion, helping maintain clarity in large-scale databases.
4. Enhanced data governance: IDs support clear ownership and governance policies, as you can assign responsibility for specific records to the appropriate teams or systems.

Overall, using IDs is a fundamental practice that not only enables reliable data integration and consistency but also improves the overall manageability and value of your data.

Column: Data Cleansing

Data cleansing is the process of identifying, correcting, or removing errors, inconsistencies, and inaccuracies in datasets to improve data quality and reliability. It ensures that data is complete, accurate, and consistent, making it more suitable for analysis, reporting, and decision-making. Data cleansing involves addressing issues such as missing values, duplicate records, incorrect formatting, and outliers. This process is crucial for maintaining the integrity of data, especially when integrating multiple datasets or preparing data for advanced analytics.

- Data Cleansing Techniques
 - Removing duplicates: Identifying and eliminating duplicate records to avoid redundant or misleading data
 - Handling missing data: Filling in missing values using techniques like interpolation, imputation, or deletion where necessary
 - Standardizing data: Ensuring data follows consistent formats, such as standardizing date formats or capitalization
 - Validating data: Checking data against predefined rules or reference datasets to ensure accuracy and consistency
 - Correcting errors: Identifying and rectifying typos, incorrect entries, or mismatched values

Effective data cleansing ensures high-quality data, which leads to better insights and more accurate decision-making.

Column: Data Matching

Data matching is the process of comparing and identifying records from different datasets—or within the same dataset—to determine whether they represent the same entity. It is widely used in scenarios such as deduplication, linking customer information across systems, fraud detection, and data integration. Data matching involves evaluating attributes such as names, addresses, phone numbers, or other identifiers to detect matches, even when there are slight variations or inconsistencies in the data. The level of matching is sometimes expressed as a score.

By consolidating duplicate or related records, data matching improves data quality and ensures accurate analysis and reporting.

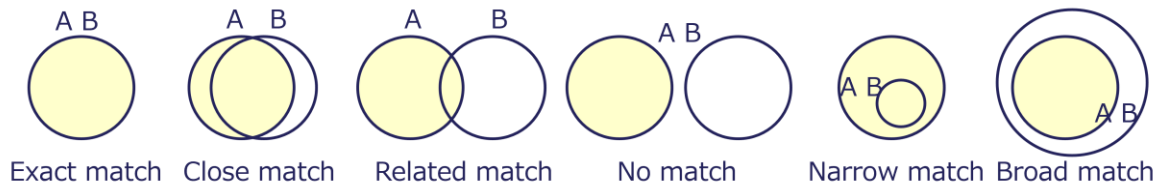


Figure 19. Data Matching Pattern

Column: Data Completion

Concepts such as data augmentation and data enrichment are used for data completion.

Data augmentation is a technique used to enhance the size and diversity of a dataset by applying various transformations to existing data. These transformations can include, for example, rotating, flipping, cropping, adjusting brightness and contrast, or adding noise for image data, paraphrasing, or back-translation for text data, or other transformations for different types of data. The goal is to increase the variety of training data, which helps improve the performance and robustness of machine learning models.

Data enrichment is the process of enhancing existing data by supplementing it with additional information from external or internal sources. This process adds value to raw data, making it more comprehensive, accurate, and useful for analysis or decision-making. Enrichment often includes adding missing data points, improving data accuracy, or integrating contextual information such as demographic, geographic, or behavioral data.

Column: Synthetic Data

While synthetic data offers expected benefits, it also carries potential risks.

Benefits

The use of synthetic data generated by AI offers significant benefits and valuable opportunities.

1. Enhancing data privacy
 - Synthetic data enables analysis and model training without exposing real personal information, helping to protect individual privacy.
2. Expanding data availability and safe testing
 - Synthetic data can fill gaps where real data are scarce, sensitive, or difficult to collect, and enables safe testing and validation of AI systems under controlled conditions.
3. Accelerating AI development
 - By generating diverse and high-quality training data, synthetic data helps improve AI performance, robustness, and fairness.
4. Reducing costs and risks
 - Using synthetic data lowers the cost of data collection and annotation, while minimizing regulatory and security risks associated with handling real data.
5. Promoting innovation and collaboration
 - Openly shared synthetic datasets can encourage cross-industry collaboration and faster research progress while maintaining confidentiality.

Risks

The synthetic content generated by AI has a significant impact.

1. Spread of fake news and misinformation
 - Synthetic content can generate highly realistic images, audio, and text, which may be intentionally misused to spread false information.
2. Invasion of privacy
 - Technologies like deepfakes can create content that mimics an individual's face or voice, potentially violating his or her privacy.
3. Damage to brand or corporate reputation
 - Maliciously created synthetic content can tarnish a company or brand reputation.
4. Legal risks and ethical challenges

- The creation and distribution of synthetic content can lead to legal and ethical issues such as copyright infringement, violation of publicity rights, or discriminatory content.
5. Decline in trust
 - An abundance of synthetic content makes it harder for consumers and businesses to trust digital content.
 6. National security and political risks
 - Synthetic content might be politically exploited.
-

Column: Content Transparency Tech

Content Transparency Technologies refer to tools, frameworks, or systems designed to provide clarity, authenticity, and accountability in digital content. These technologies aim to ensure that users can understand the origin, context, and reliability of the information they encounter. They are often used in media, advertising, social platforms, and data-driven industries to address issues such as misinformation, biased content, or non-transparent algorithms.

Key features of Content Transparency Technologies include:

- Source verification: Ensuring content authenticity by verifying the origin or creator of the material
 - Traceability: Providing a clear history of edits, ownership, and dissemination of the content
 - Content labeling: Adding metadata or markers to indicate whether content is sponsored, user-generated, or machine-generated
-

Column: Indication of Data Quality

The level of data quality indicators should naturally differ between internal management and public use. Data quality is often managed in detail within organizations by data producers and users. However, when data circulates more broadly in society, **simpler and easier-to-understand indicators** become necessary. This is similar to how home appliances are handled: detailed specifications and quality checks are used in manufacturing and logistics, but consumers only see simple indicators such as efficiency labels or safety marks at stores. Likewise, in data quality management, it is important to design different levels of indicators—detailed ones for internal management and simple ones for external sharing or trading—according to the purpose and stakeholders involved.



Figure 20. Context Shift of Data Quality Indicators for Internal Management and External Use

4.1.5 Data Processing

4.1.5.1 General

- Description
 - Data processing creates service value by processing data effectively and efficiently.
 - This process includes identifying root causes within the data path and providing actionable feedback when errors occur.
 - The workflow emphasizes maintaining data integrity and seamless operation throughout the processing pipeline.
- Issue
 - Data inconsistencies, mismatches, or incomplete records that disrupt the processing pipeline
 - Difficulty in maintaining reliability and accuracy in processed data
 - Time-consuming root cause identification without robust diagnostic tools or clear process models
- Requirement
 - Comprehensive processes and robust data models for identifying and addressing the root causes of errors
 - Monitoring and diagnostic mechanisms, including real-time monitoring where necessary, for quicker error resolution and minimized service disruptions
- Value
 - Enhanced accuracy and trustworthiness of services through effective error handling and high-quality data processing
 - Improved customer confidence, reduced operational risks, and support for data-driven innovation

Process the data that has been entered. As there is a possibility of data inconsistencies occurring during the processing, a feedback mechanism is required.

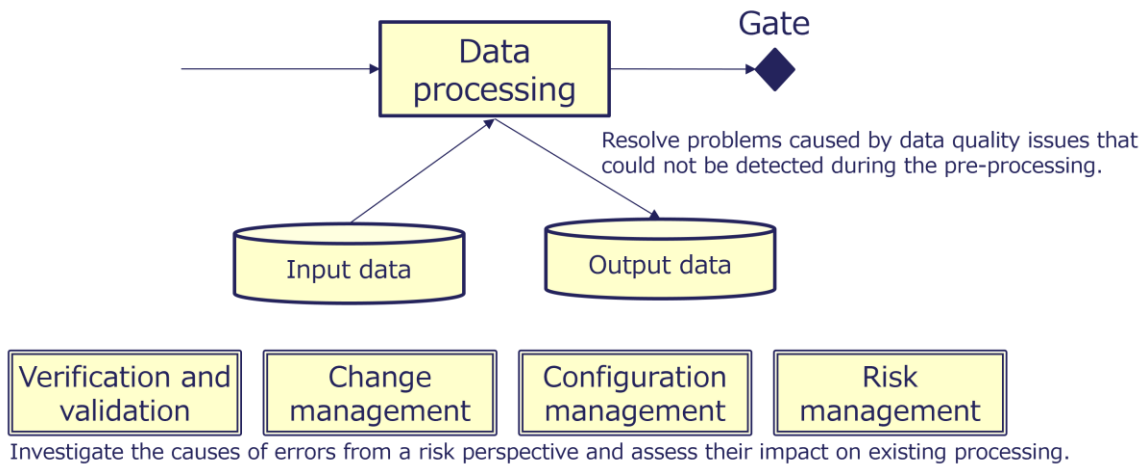


Figure 21. Data Processing Process

4.1.5.2 Data Processing

- Procedure
 1. Check consistency.
 2. Process data.
 3. Report errors in processing.
- Checkpoint
 - Is an error notification provided when there is an anomaly in the data?
 - Is the data processing visualized for verification?

4.1.6 AI System

4.1.6.1 General

- Description
 - The AI system process covers training AI models using curated datasets.
 - It includes implementation, utilization, and operation of trained AI models.
 - It includes continuous retraining and refinement based on output evaluations.
 - For example, RAG techniques can be used to enhance the accuracy and reliability by referencing external knowledge sources.
- Issue
 - Risk of hallucination, including false or inaccurate outputs when data is limited or unclear

- Risk of societal and/or data bias in outputs, raising issues of fairness, inclusion, and ethics
- High time, resource, and oversight requirements for training and relearning, reducing scalability
- Requirement
 - Processes for improving the accuracy and reliability of answers to user queries
 - Mechanisms for providing supporting reference information, such as verifiable sources or reasoning, to build trust and ensure transparency in the decision-making process
- Value
 - Higher service levels through accurate, fast, and contextually relevant responses
 - Increased trust in the AI system among users and stakeholders through reliable, unbiased, and transparent outputs

Using high-quality, fit-for-purpose data improves accuracy during operation.

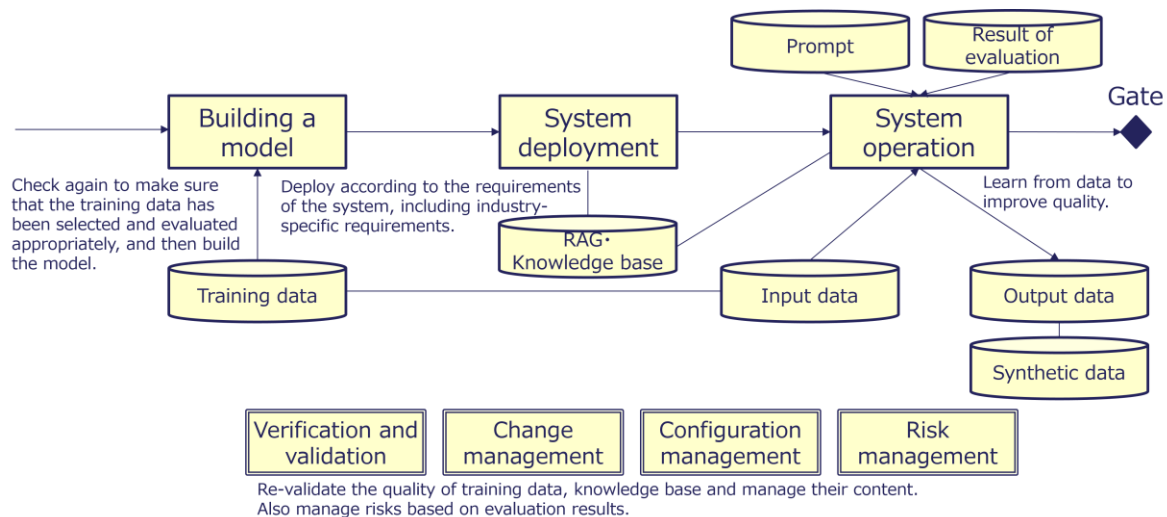


Figure 22. AI System Process

4.1.6.2 Building a Model

- Procedure
 1. Check the required level for AI.
 2. Decide on the training data.
 - Prevent biased data and inappropriate data.
 3. Decide whether to use RAG.

4. Check whether there is any personal or intellectual property information.
 5. Train the model.
 6. Test the model.
- Checkpoint
 - Do you use reliable data for your training data?
 - Where appropriate, have you implemented RAG to improve accuracy?
 - If there is information that could lead to personal data or intellectual property, do you take measures to prevent this information from being displayed, for example by excluding the unnecessary information from the training data?

4.1.6.3 System Deployment

- Procedure
 1. Deploy the system.
 2. Train users.
- Checkpoint
 - Do you have a deployment policy and follow it?

4.1.6.4 System Operation

- Procedure
 1. Operate the system.
 2. Continuously monitor.
- Checkpoint
 - Does reusing the output introduce any bias into the AI system?

Column: Bias

Reducing bias is essential for implementing AI. By addressing bias, AI systems can foster broader societal acceptance, ensure ethical use, and maintain long-term reliability and sustainability.

1. Preventing unfairness and discrimination
 - AI trained on biased data may make decisions that disadvantage specific groups based on gender, age, race, or location, leading to ethical concerns.
2. Ensuring trust and fairness

- Bias undermines trust in AI systems. To create trustworthy AI solutions, fairness and impartiality must be maintained.
3. Avoiding business risks
 - Bias-related controversies can harm a company's reputation and expose it to legal risks.
 4. Minimizing societal impact
 - The widespread adoption of AI amplifies its societal influence. If biased, AI can exacerbate inequalities or deepen societal divides.
 5. Compliance with legal and ethical standards
 - Increasing regulations around AI emphasize fairness and bias mitigation. Non-compliance can lead to penalties or exclusion from specific markets.
-

Column: Model Collapse

What is Model Collapse?

A phenomenon where generative AI models trained repeatedly on synthetic data (rather than real human data) gradually lose diversity and accuracy over time.

The following figure shows how model collapse progresses through the repeated use of synthetic data.

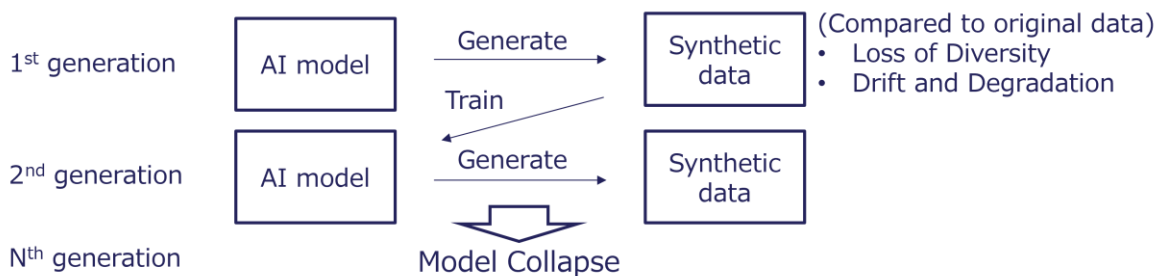


Figure 23. Mechanism of Model Collapse

Prevention and Mitigation

- Ensure clear distinction and proper governance between real and synthetic data.
- Periodically retrain using verified real-world datasets.
- Implement continuous quality monitoring and feedback loops for dataset integrity.

Source: Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, Yarin Gal, 2024, “[AI models collapse when trained on recursively generated data](#)”

4.1.7 Evaluation of Output

4.1.7.1 General

- Description
 - Evaluation of output serves as the final verification before AI system outputs are provided to end users.
 - This process includes assessing outputs for ethical considerations, information accuracy, potential biases, and alignment with intended purposes.
 - It uses automated systems and human reviewers to evaluate outputs against reliable information sources and pre-established evaluation datasets.
 - It includes reporting identified errors back to data holders or system owners to enable continuous improvement.
- Issue
 - Risk of inaccurate or inappropriate responses that reduce user trust and system reliability
 - Deviations from authoritative definitions or standards causing inconsistencies
 - Insufficient mechanisms to detect, categorize, and correct errors
- Requirement
 - Robust measures for addressing accuracy, ethics, and bias to prevent inaccurate or inappropriate outputs
 - Clear flagging of outputs that require user confirmation
 - Feedback loops for preventing repeated errors and improving system performance
- Value
 - Higher service levels through outputs that meet or exceed user expectations for reliability, accuracy, and usability
 - Increased trust and credibility among users, stakeholders, and society through ethical and transparent operation of the AI system

Evaluation of outputs is important as a safeguard for the system. Prevent inappropriate responses.

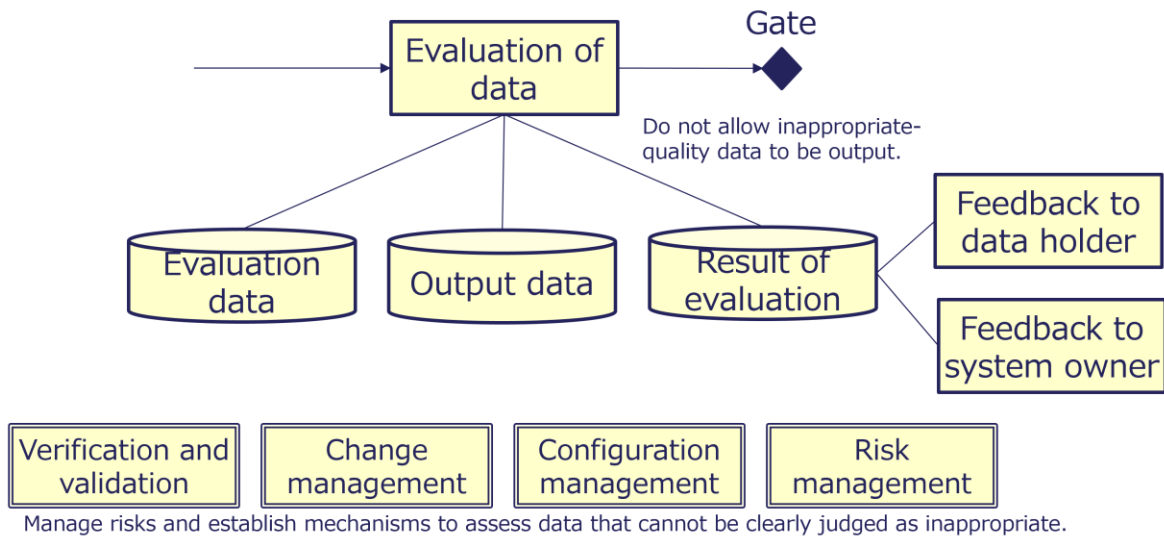


Figure 24. Evaluation of Output Process

4.1.7.2 Evaluation of Output

- Procedure
 1. Prepare the evaluation data.
 2. Verify outputs against reliable data.
 3. Remove inappropriate responses based on the evaluation data.
 4. Create error reports.
 5. Provide feedback to data holders.
 6. Provide feedback to system owners.
- Checkpoint
 - Are safeguards in place to prevent unethical responses?
 - If a decision needs to be made on a response, does the process involve human intervention?
 - Do you have mechanisms in place to investigate the causes of inappropriate responses and provide feedback?
 - Are there any outputs that do not match the reliable data?

Column: ASoT

The **Authoritative Source of Truth (ASoT)** refers to the definitive, trusted source of information for a specific system, process, or context. It is the location where accurate, up-to-date, and complete data is maintained, ensuring consistency and

reliability across systems or teams that rely on that information. By verifying data using ASoT, you can eliminate incorrect data and improve data quality. Open data provided by the government serves as essential social infrastructure for enhancing data quality.



Figure 25. Illustration of Data Verification Using ASoT and Open Data

4.1.8 Delivering Results

4.1.8.1 General

- Description
 - Delivering results covers providing processed data and AI outputs to end users in a clear and accessible manner.
 - This process includes presentation measures that prevent misunderstandings and support trust and usability.
 - It includes metadata management to provide context and reliability.
 - It may include watermarking to secure data and protect intellectual property.
- Issue
 - Risk of misunderstandings due to poorly presented data
 - Interoperability barriers caused by non-machine-readable or non-standardized interfaces
 - Risk of inappropriate use of outputs without proper guidelines or safeguards
- Requirement
 - Presentation processes that support ease of understanding and accessibility for diverse audiences
 - Processes for preserving data sovereignty, ownership, control, and regulatory compliance
 - Standardized formats and interfaces for interoperability across systems
- Value
 - Greater societal benefits of AI through clear data sharing, informed decisions, and trust

- Support for ethical use and sustainable growth in data-driven ecosystems

Outputs should be made publicly available or provided for use in other systems. It is important to prevent misinterpretation and its potential social impact.

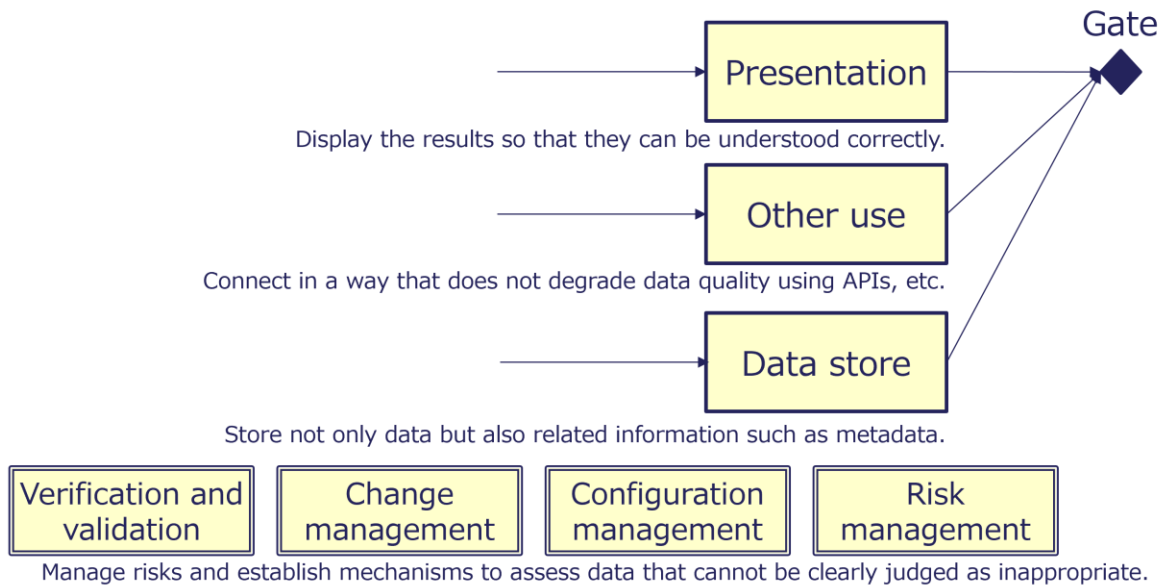


Figure 26. Delivering Results Process

4.1.8.2 Presentation

- Procedure
 1. Manage access controls.
 2. Present the processing results.
 - Easy to understand
 - Visualization
 3. Gather the user's opinion.
- Checkpoint
 - Are the answers and expressions misleading?
 - Is there a mechanism to check the basis of the answers, for example, by making the provenance information available for checking?

4.1.8.3 Other Use

- Procedure
 1. Manage access controls.

2. Provide APIs and relevant information.
- Checkpoint
 - Does it provide a machine-readable interface?

4.1.8.4 Data Store

- Procedure
 1. Store the processing results.
 2. Back up the data.
 3. Optimize data storage.
 - Checkpoint
 - Are protective measures taken to ensure that data is not lost?
-
-

Column: Social and Human Impact

Low-quality data may cause a system to provide incorrect information, which may harm, mislead, or disadvantage users. Furthermore, such incorrect information may lead to accidents or broader economic impacts. In addition to evaluating data before it is provided, it is also important to prepare countermeasures in case any negative impact occurs. If incorrect information is released, it is important to have reliable sources that can provide and verify accurate data.

4.1.9 Decommissioning

4.1.9.1 General

- Description
 - Decommissioning covers the formal completion of the data quality control process and the official retirement of data.
 - This process includes measures for clearly marking data as decommissioned to prevent accidental or unintended use.
 - It supports the integrity of AI systems by preventing errors caused by outdated data.
- Issue
 - Risk of users unknowingly using decommissioned or outdated data and reaching inappropriate conclusions or results
- Requirement
 - Systems and processes for prominently notifying users that data has been decommissioned

- Metadata updates, automated alerts, and clear documentation to eliminate ambiguity about the data's status
- Value
 - Reduced errors and rework caused by the use of outdated data
 - Improved AI system reliability, trust, operational efficiency, and data governance through clear communication of decommissioning status

Stop providing data. Provide advance notice of decommissioning and indicate that the data is no longer maintained or guaranteed thereafter.

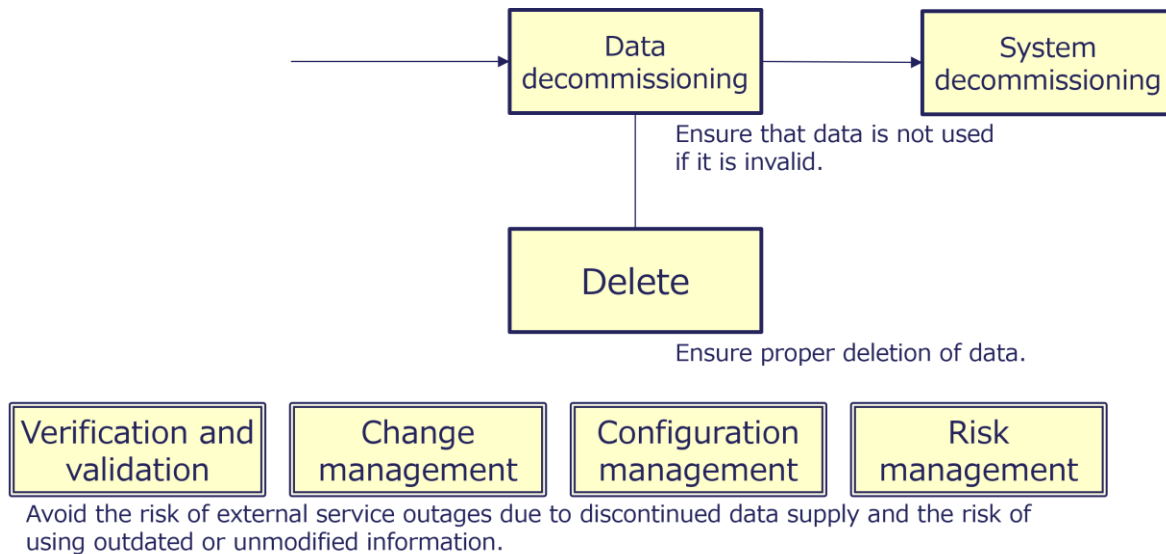


Figure 27. Decommissioning Process

4.1.9.2 Data Decommissioning

- Procedure
 1. Notify users of the decommissioning of the data.
 2. Provide information on data that has been decommissioned.
 3. Transfer not only the data itself but also its metadata and related documents.
- Checkpoint
 - Did you give sufficient notice before suspending the data?
 - If you are transferring to another party, did you provide sufficient information for them to take over?

4.1.9.3 Delete

- Procedure

1. If necessary, archive the data.
 2. Erase the data so that it cannot be restored.
 3. When deleting parts, give prior notice.
- Checkpoint
 - Did you check the results of the deletion?

4.1.9.4 System Decommissioning

- Procedure
 1. Dispose of system components.
- Checkpoint
 - Are system components properly decommissioned?

4.1.10 Processes throughout the Life Cycle

4.1.10.1 General

- Description
 - Processes throughout the data life cycle cover verification and validation, change management, configuration management, and risk management.
 - These initiatives support the improvement of data quality across processes.
- Issue
 - Lack of management for the system or service as a whole, making review and rapid response difficult when problems arise
- Requirement
 - Management of data quality-related information in a system that is easy for stakeholders to understand
 - Traceability of data quality-related information
- Value
 - Improved compliance and transparency
 - Increased trust through improved transparency

4.1.10.2 Verification and Validation

- Procedure
 1. Define the characteristics, targets and tolerance levels required for the system's purpose.
 2. Manage the system throughout the life cycle.
- Checkpoint

- Is data quality control placing a burden on operations?

4.1.10.3 Change Management

- Procedure
 1. Define a basic policy for changes, such as data integration, processing, and modification.
 2. Define policies for data degradation over time, e.g., data refreshing.
 3. Record changes.
- Checkpoint
 - Is the change history readily available?

4.1.10.4 Configuration Management

- Procedure
 1. Manage software configurations.
 2. Manage data configurations.
 3. Manage the configurations of relevant documents.
- Checkpoint
 - Do you manage the list of software, data and documents?

4.1.10.5 Risk Management

- Procedure
 1. Define data quality risks and response policies.
 2. Manage data quality risks using access control and continuous monitoring.
 3. Ensure that essential risk factors are addressed.
 4. Develop a Business Continuity Plan (BCP).
- Checkpoint
 - Do you understand the risk of your system or service?
 - When a risk is discovered, does your organization have a culture of reviewing the situation from the root cause?

Column: Access Control

Access control is essential for maintaining data quality by preventing data poisoning and other forms of unauthorized manipulation. By defining and regulating which users, devices, or processes are permitted to access, modify, or delete data, organizations can significantly reduce the risk of malicious attacks and accidental errors. Key

elements of effective access control include implementation of clearly defined user roles and privileges, enforcing multi-factor authentication, regularly auditing access logs, and continuously monitoring for unusual or unauthorized activities. This structured approach ensures data remains accurate, consistent, and secure throughout its entire life cycle.

Column: Improving Data Quality Using AI

AI and data quality have a mutually reinforcing relationship. While data quality management enhances AI performance, AI can also be leveraged to improve data quality.

AI-assisted data quality management involves using artificial intelligence to detect, correct, and prevent data errors. AI algorithms can automatically identify inconsistencies, missing values, and duplicates by learning from data patterns. They also suggest improvements, monitor data quality over time, and adapt to changing data environments. This helps ensure accurate, reliable, and up-to-date information for decision-making.

The following figure illustrates that AI and data quality have a bidirectional, mutually reinforcing relationship rather than a one-way relationship.

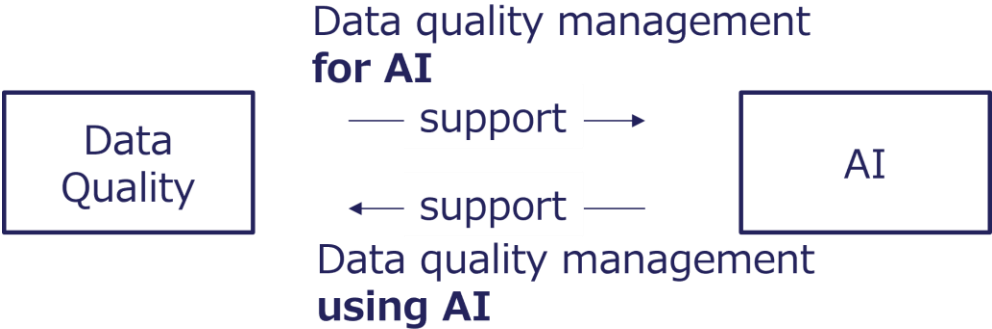


Figure 28. Mutual Relationship Between AI and Data Quality

4.2 Governance Cycle View

4.2.1 General

Managing data quality involves performing individual activities correctly and establishing organizational controls to ensure that those activities are sustainable and operate smoothly.

This view focuses on the governance of data quality, ensuring that the processes defined in the Process View are appropriately controlled, monitored, and improved at an organizational level. In this context, “processes” refer to management processes that oversee and support the execution of data life cycle processes.

In some cases, management processes may overlap with life cycle processes. This reflects the fact that certain processes, such as validation and risk management, have both operational and governance aspects.

This view aligns the data quality governance structure with ISO 8000-61.

The figure below shows the basic structure of the data quality management process defined in ISO 8000-61. To achieve continuous improvement, this process follows the PDCA cycle (in this case, Plan–Control–Assurance–Improvement). It also illustrates the subprocesses within each phase. Data-related supporting processes provide information and technologies that enable implementation.

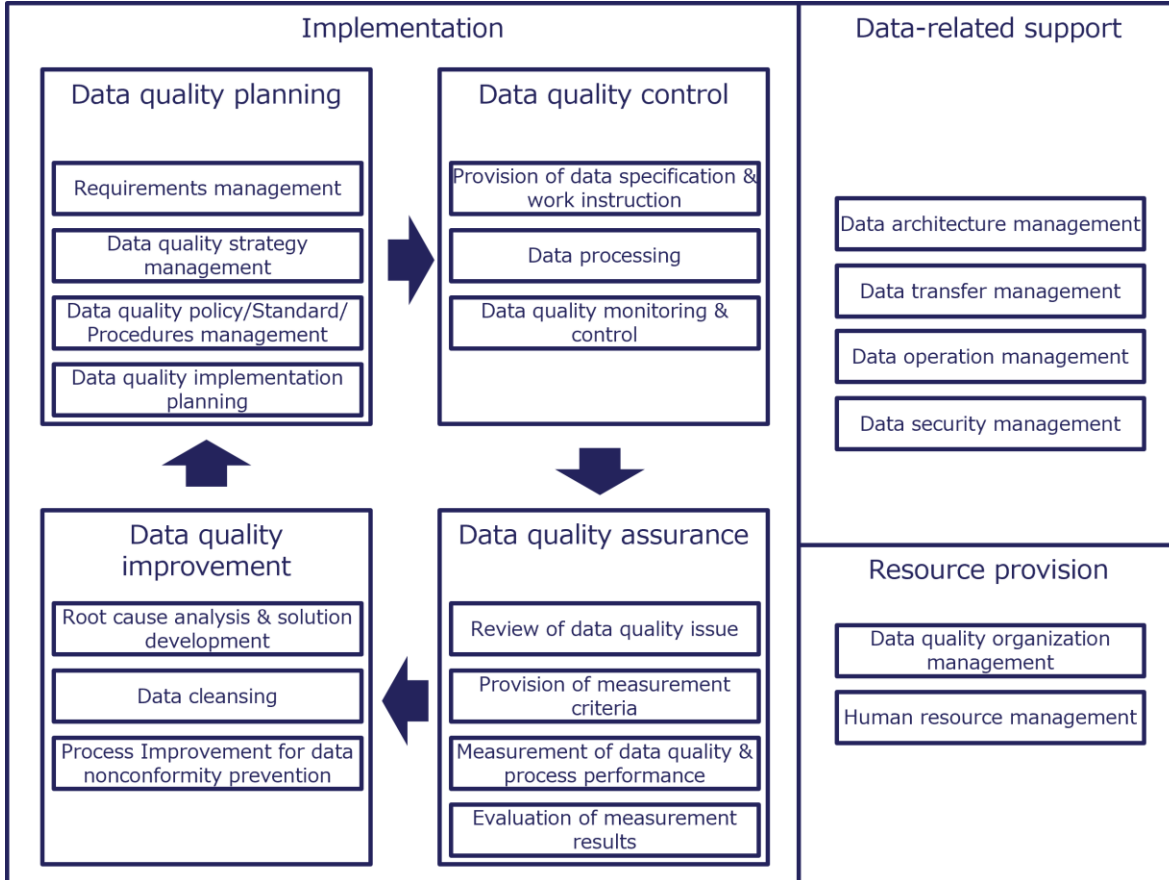


Figure 29. Governance Cycle for Data Quality Management

4.2.1.1 Structure of This Section

This section uses a two-part structure for each management process.

- General
 - Provides an overview of each management process, including its purpose, role, and value in ensuring data quality, along with relevant risks, scope, stakeholders, and practical considerations.
- Process
 - Presents subprocesses carried out within each management process, along with representative activities that may be performed in association with them.

4.2.2 Data Quality Planning

4.2.2.1 General

Data quality planning defines the organizational approach for managing data quality across the AI life cycle. It supports accurate, reliable, and consistent data for AI systems, reducing risks such as bias, errors, and poor decision-making. This process includes identifying data requirements, defining quality standards, establishing policies, and developing actionable plans. A key practical consideration is ensuring that executives understand the importance of data quality.

4.2.2.2 Requirements Management

- Business needs identification: Determine the specific objectives of the AI system and the data needed to achieve them.
- Data quality dimensions definition: Establish the critical dimensions (e.g., accuracy, completeness, consistency, timeliness) relevant to the project.
- Current data assessment: Conduct a gap analysis of existing data to identify areas needing improvement.
- Requirements documentation: Clearly document the data quality requirements to guide subsequent processes.

4.2.2.3 Data Quality Strategy Management

- Vision development: Define the long-term goals for data quality aligned with the AI system's objectives.
- Measurable target setting: Establish key performance indicators (KPIs) to evaluate data quality over time.
- Stakeholder engagement: Ensure alignment with all stakeholders, including data engineers, analysts, and decision-makers.
- Risk management: Identify and mitigate potential risks associated with low-quality data in AI applications.

4.2.2.4 Data Quality Policy/Standard/Procedure Management

- Policy definition: Establish organizational policies for data governance and quality management.
- Standard creation: Develop specific standards for data collection, storage, processing, and validation to ensure consistency.
- Procedure documentation: Outline detailed procedures for maintaining and improving data quality, such as periodic audits and validation protocols.
- Compliance checking: Ensure that all policies and standards align with relevant legal and regulatory requirements.

4.2.2.5 Data Quality Implementation Planning

- Action plan development: Create a step-by-step plan for implementing data quality measures, including timelines and responsibilities.
- Role assignment: Define clear roles and responsibilities for team members involved in data quality management.
- Tool and process implementation: Deploy tools for data cleansing, validation, monitoring, and reporting.
- Monitoring and refinement: Continuously monitor data quality and refine the implementation plan as needed based on feedback and performance metrics.

4.2.3 Data Quality Control

4.2.3.1 General

Data quality control aims to ensure that data used in AI systems meets required quality levels for accuracy, consistency, and reliability. It supports AI performance optimization, bias reduction, and trustworthy results. This process includes clear specifications, robust processing techniques, and continuous monitoring to identify and address quality issues. A key practical consideration is avoiding excessive burden on the workplace by using technical measures such as automated checks.

4.2.3.2 Provision of Data Specifications and Work Instructions

- Data requirements definition: Define clear data requirements, including formats, structures, and acceptable value ranges.
- Work instruction provision: Provide detailed instructions for data collection, labeling, and processing to ensure uniformity.
- Metadata guideline establishment: Establish guidelines for metadata creation to track data sources, timestamps, and context information.
- Validation checklist development: Develop a validation checklist for contributors to ensure adherence to data quality standards.

4.2.3.3 Data Processing

- Data cleaning and preprocessing: Clean and preprocess data to remove inconsistencies, duplicates, and irrelevant information.
- Data normalization: Normalize data to ensure compatibility across systems (e.g., format conversions, scaling).
- Annotation and labeling: Annotate and label data accurately for supervised learning models.
- Automated validation implementation: Implement automated data validation tools to identify anomalies and errors during processing.

- Processing documentation: Document each processing step to maintain traceability and reproducibility.

4.2.3.4 Data Quality Monitoring and Control

- Quality metrics establishment: Establish key quality metrics (e.g., accuracy, completeness, consistency, timeliness).
- Monitoring implementation: Implement timely monitoring systems, including real-time monitoring where necessary, to detect data quality issues promptly.
- Regular audit scheduling: Schedule regular audits to review and validate data integrity.
- AI tool use: Use AI tools to identify patterns of errors or potential biases in datasets.
- Feedback loop creation: Create feedback loops to continuously improve data quality standards based on findings.

4.2.4 Data Quality Assurance

4.2.4.1 General

Data quality assurance (DQA) provides assurance that data used in AI development and operations meets required standards of accuracy, consistency, completeness, and reliability. It helps prevent errors, biases, and inefficiencies in AI systems. This process includes identifying potential data quality issues, defining evaluation criteria, measuring data quality, and analyzing results to guide continuous improvements.

4.2.4.2 Review of Data Quality Issues

- Potential issue identification: Identify potential issues such as missing values, inconsistencies, inaccuracies, or redundancies in the dataset.
- Root cause analysis: Analyze root causes of data quality problems, including errors in data collection, processing, or storage.
- Issue documentation: Document known issues and assess their potential impact on AI model performance.

4.2.4.3 Provision of Measurement Criteria

- Measurement criteria definition: Define clear and measurable criteria for data quality attributes (e.g., accuracy, completeness, timeliness, consistency, and validity).
- Benchmark and threshold establishment: Establish benchmarks and thresholds that data are expected to meet to be considered suitable for use.
- Criteria alignment: Align criteria with the specific requirements of AI models and business objectives.

4.2.4.4 Measurement of Data Quality and Process Performance

- Dataset evaluation: Conduct systematic evaluations of datasets against the established measurement criteria.
- Data profiling: Use data profiling tools to detect anomalies and assess quality attributes.
- Workflow monitoring: Monitor data processing workflows to ensure consistent adherence to quality standards.

4.2.4.5 Evaluation of Measurement Results

- Measurement outcome analysis: Analyze measurement outcomes to identify gaps and areas for improvement.
- Impact quantification: Quantify the impact of data quality issues on AI system performance and decision-making.
- Reporting and dashboard creation: Generate reports and dashboards to communicate findings to stakeholders.

4.2.5 Data Quality Improvement

4.2.5.1 General

Data quality improvement is a management process for improving the reliability and effectiveness of AI systems by addressing data quality problems. It supports sustainable quality by identifying and resolving data issues, improving data consistency, and implementing preventive measures against data deterioration over time. This process includes root cause analysis, solution development, data cleansing, and process improvement to prevent data nonconformity.

Rather than merely modifying the data, it is essential to address root causes such as difficulties in data entry or inefficiencies in the process itself, to achieve sustainable improvement.

4.2.5.2 Root Cause Analysis and Solution Development

- Issue source identification: Investigate the underlying causes of data quality problems, such as incorrect data entry, system errors, or outdated data.
- Targeted solution development: Design and implement specific corrective actions to address identified root causes, such as updating workflows, improving validation rules, or automating data entry processes.
- Result monitoring and validation: Continuously evaluate the effectiveness of implemented solutions to ensure the issues are resolved and do not reoccur.

4.2.5.3 Data Cleansing

- Inaccuracy identification: Detect duplicate, incomplete, or inconsistent records within the dataset.

- Data format standardization: Ensure uniformity in data presentation, such as consistent date formats, unit measurements, and naming conventions.
- Error correction or removal: Modify inaccurate entries, fill missing values, or remove irrelevant or outdated data to maintain dataset integrity.
- Cleansing process automation: Leverage tools and algorithms to automate repetitive data cleansing tasks and reduce manual errors.

4.2.5.4 Process Improvement for Data Nonconformity Prevention

- Data governance policy establishment: Define and enforce clear policies for data collection, management, and usage to prevent quality issues.
- Data validation mechanism enhancement: Implement validation checks, including real-time validation where necessary, during data entry or ingestion to identify nonconformities early.
- Stakeholder training: Educate employees and data handlers about best practices for maintaining data quality and the importance of adhering to standards.
- Data process monitoring and auditing: Regularly assess data workflows to detect potential sources of nonconformity and ensure compliance with established policies.

4.2.6 Data-related Support

4.2.6.1 General

Data-related support provides the information, technologies, and operational support needed for effective data quality governance in AI development and utilization. It supports the integrity, consistency, and reliability of data across its life cycle. This process includes data architecture management, data transfer management, data operations management, and data security management. These efforts establish frameworks, processes, and tools that optimize data management while addressing compliance, security, and operational needs.

4.2.6.2 Data Architecture Management

- Data architecture and metadata repository management: Define data models, schemas, and standards to ensure consistency, and develop and manage metadata repositories to enhance data discoverability and traceability.
- Data lineage tracking: Implement data lineage tracking to monitor the flow and transformation of data.
- Scalability and flexibility assurance: Ensure scalability and flexibility to accommodate evolving AI demands.
- Architecture alignment: Align architecture with organizational goals and compliance requirements.

4.2.6.3 Data Transfer Management

- Secure data transfer protocol establishment: Establish protocols for secure data transfer (e.g., encryption, VPNs).
- Data flow monitoring and audit maintenance: Monitor and optimize data flow to prevent bottlenecks and ensure timely availability, including real-time availability where required, and maintain logs and audits of data transfers for accountability and traceability.
- Cross-border data transfer policy definition: Define policies for cross-border data transfers to comply with legal regulations (e.g., GDPR, CCPA).
- Data transfer automation: Automate data transfer processes where applicable to reduce manual errors.

4.2.6.4 Data Operations Management

- Data validation and cleansing: Conduct regular data validation and cleansing to maintain accuracy.
- Data storage management: Manage data storage solutions to ensure accessibility and scalability.
- Data workflow establishment and performance monitoring: Establish workflows for data ingestion, processing, and integration, monitor system performance, and resolve issues related to data operations.
- Dataset version control: Implement version control for datasets to track changes and ensure consistency.

4.2.6.5 Data Security Management

- Access control implementation: Implement access controls, including Role-Based Access Control (RBAC) and multifactor authentication.
- Security audit and vulnerability assessment: Conduct regular security audits and vulnerability assessments.
- Encryption and incident monitoring: Deploy encryption protocols for data at rest and in transit, monitor for suspicious activities, and respond to security incidents promptly.
- Regulatory compliance assurance: Ensure compliance with global and local data protection regulations.

4.2.7 Resource Provision

4.2.7.1 General

Resource provision ensures that the organizational structures, human resources, and operational capabilities needed for effective data quality governance are in place. It supports the provision of high-quality data for AI systems while minimizing risks and improving the accuracy, reliability, and fairness of AI models. This process includes

organizing teams, allocating skilled personnel, and defining clear roles and responsibilities.

4.2.7.2 Data Quality Organization Management

- Data quality governance team establishment: Establish a dedicated Data Quality Governance team or council.
- Role and responsibility definition: Define roles and responsibilities for data quality management (e.g., Data Stewards, Data Quality Analysts).
- Reporting line and accountability setup: Set up clear reporting lines and accountability for data quality issues.
- Cross-functional collaboration framework creation: Create cross-functional collaboration frameworks to involve stakeholders from various departments.
- Organizational structure review and update: Regularly review and update organizational structures to adapt to evolving data and AI needs.

4.2.7.3 Human Resource Management

- Data quality professional recruitment: Recruit data quality professionals with expertise in data governance, data architecture, and AI integration.
- Training program provision: Provide training programs to enhance employees' skills in data quality assessment, cleansing, and validation.
- Career path and development plan definition: Define career paths and development plans for data governance professionals to improve retention.
- Staffing level assurance: Ensure adequate staffing levels for ongoing data quality monitoring and improvement tasks.
- Data quality culture promotion: Promote a culture of data quality awareness and accountability across the organization.

Column: Human Resource and Team

The data governance team will work with the AI team to improve data quality.

The figure below illustrates the key roles involved in collaboration between the AI team (CAIO and AI Office) and the data governance team (CDO and Data Steward). Since AI and data are interdependent, the people responsible for them must also collaborate closely.

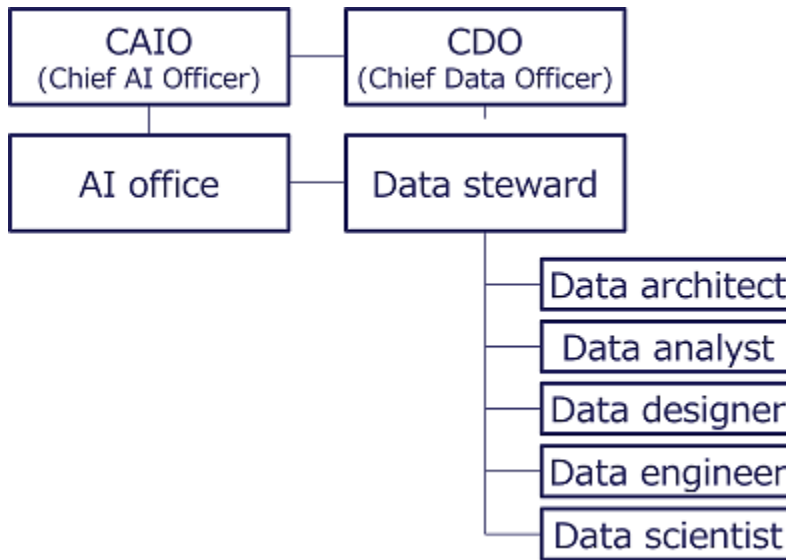


Figure 30. Collaboration Between AI and Data Governance Teams

Business fields using the data also need literacy training, including data quality.

The following diagram illustrates a cyclical process for human resource development and the enhancement of literacy. Learning content, as a systematically organized body of knowledge, supports knowledge acquisition. That knowledge is then applied through use cases, leading to practical experience. These experiences are shared within the community, encouraging further action in the field. Actions taken in the field contribute to the creation and improvement of new learning content, thereby completing the cycle once. Skill standards function as the foundation that underpins this entire process, providing a common framework that ensures consistency and supports continuous development at every stage.



Figure 31. Human Resource Development Cycle

4.3 Gateway View (Characteristics)

4.3.1 General

To maintain data quality, it should be checked at process boundaries and transactions. In the case of sensor data, this may be a periodic inspection. During the inspection process, each characteristic should be checked, assuming the type of data and events that may degrade the quality of the data. The checking of characteristics should be automated where possible, so as not to increase the workload on the site.

4.3.1.1 Structure of This Section

ISO/IEC 25012 and ISO/IEC 5259-2 classify data quality characteristics based on whether they relate to data content, system aspects, both, or additional characteristics specific to the context of AI. The following sections describe each characteristic categorized according to the above classification, with evaluation points and inappropriate examples. In addition, this document addresses data quality characteristics relevant to sensor data.

- Evaluation Points
 - Describe examples of evaluation perspectives that may be used to assess whether each characteristic has been met.

- Inappropriate Examples
 - Provide representative examples of situations where the characteristic is not appropriately satisfied.

4.3.2 Inherent Data Quality Characteristics

4.3.2.1 Accuracy

Accuracy ensures that the data reflect real-world values correctly. High accuracy is vital for AI to produce reliable outcomes, avoiding false predictions or flawed decisions.

- Evaluation Points
 - Data match verified sources
 - Error rate is within acceptable thresholds
 - Outliers are justified or corrected
- Inappropriate Examples
 - Incorrect labels in training datasets
 - Mismatched units in numerical data
 - Spelling errors in text data

4.3.2.2 Completeness

Completeness measures whether all required data points are present. Missing values can distort AI predictions and compromise system effectiveness.

- Evaluation Points
 - All critical fields are filled
 - Null or missing rates are minimal
 - Completeness checks align with AI model requirements
- Inappropriate Examples
 - Missing demographic attributes in user profiles
 - Partial transaction records
 - Null values in key input variables

4.3.2.3 Consistency

Consistency ensures that data values are uniform across datasets and time. Inconsistencies can cause AI to misinterpret patterns.

- Evaluation Points
 - Formatting and units are uniform
 - Datasets are synchronized over time

- No conflicting entries within linked data
- Inappropriate Examples
 - Different date formats in a dataset
 - Duplicate records with conflicting details
 - Mismatched entries in integrated datasets

4.3.2.4 Credibility

Credibility measures the trustworthiness of data sources. Reliable sources improve AI model validity and trust.

- Evaluation Points
 - The data origins are verified and reputable
 - Data provenance is documented
 - Sources are peer-reviewed or authenticated
- Inappropriate Examples
 - Data from unknown/unverified sources
 - Fake or manipulated datasets
 - User-generated content without validation

4.3.2.5 Currentness

Currentness refers to whether the data is recent enough for the intended use. Stale data can result in outdated AI insights or actions.

- Evaluation Points
 - Data aligns with the latest context
 - Data is acquired with appropriate temporal coverage
 - Regular updates are maintained
- Inappropriate Examples
 - Outdated stock prices in financial models
 - Old weather data in climate predictions
 - Historical user preferences in real-time applications

4.3.3 Inherent and System-Dependent Data Quality Characteristics

4.3.3.1 Accessibility

Accessibility ensures that data is available to authorized users and systems without barriers. It involves proper storage, robust APIs, and universal design principles. Accessibility supports efficient AI model training and application by enabling smooth data flow.

- Evaluation Points
 - Data can be accessed across platforms and devices
 - APIs are well-documented and error-free
 - The data meets accessibility standards for all user demographics
- Inappropriate Examples
 - Data locked in proprietary formats
 - Missing API documentation
 - Non-compliance with Americans with Disabilities Act (ADA) accessibility standards

4.3.3.2 Compliance

Compliance ensures that data management adheres to laws, regulations, and industry standards like GDPR or CCPA. Proper compliance avoids legal risks, protects privacy, and builds trust, which is critical for AI system reliability and public acceptance.

- Evaluation Points
 - The system meets all relevant legal standards
 - Proper data consent mechanisms are implemented
 - Regular compliance audits are conducted
- Inappropriate Examples
 - Collecting data without consent
 - Ignoring jurisdictional privacy laws
 - Lacking audit trails for data handling

4.3.3.3 Confidentiality

Confidentiality ensures that sensitive data is protected from unauthorized access or breaches. Essential measures include encryption, access controls, and anonymization techniques. Maintaining confidentiality safeguards trust and reduces the risks associated with data misuse.

- Evaluation Points
 - Strong encryption is applied to data at rest and in transit
 - Role-based access control systems are adopted
 - Security protocols are regularly updated
- Inappropriate Examples
 - Storing sensitive data in plaintext
 - Sharing private data without user consent
 - Weak or outdated access controls

4.3.3.4 Efficiency

Efficiency in data management minimizes processing time and resource usage without compromising quality. Efficient data ensures faster AI training and deployment, cost savings, and scalability.

- Evaluation Points
 - The data storage structures are optimized
 - Data access latency is minimized
 - The ETL pipelines are efficient
- Inappropriate Examples
 - Redundant data processing steps
 - High latency in API calls
 - Excessive resource consumption for simple tasks

4.3.3.5 Precision

Precision ensures the degree of detail or granularity in data. The higher the precision, the more exact and detailed the information that can be obtained.

- Evaluation Points
 - Data values are recorded with sufficient detail
 - Units and scales are clearly specified
 - Temporal and spatial information is recorded with appropriate granularity
- Inappropriate Examples
 - Too coarse latitude and longitude values
 - Values with inconsistent levels of precision, making calculations difficult
 - The processing or storage burden caused by excessively high resolution

4.3.3.6 Traceability

Traceability tracks the origin, transformations, and use of data, ensuring accountability and reproducibility. Detailed logs and metadata enhance transparency, critical for debugging and compliance.

- Evaluation Points
 - Data lineage is documented in a comprehensive manner
 - Transformation logs are maintained
 - Traceability is ensured through unique IDs
- Inappropriate Examples

- Missing source details for datasets
- Unlogged data modifications
- Inconsistent versioning of datasets

4.3.3.7 Understandability

Understandability ensures that data can be interpreted correctly by both humans and machines. Clear labeling, metadata, and intuitive structures improve usability, essential for effective AI training.

- Evaluation Points
 - Data is labeled clearly and comprehensively
 - Metadata aligns with schema standards
 - The data structure is consistent and logical
- Inappropriate Examples
 - Vague or missing labels
 - Complex, undocumented structures
 - Misleading metadata annotations

4.3.4 System-Dependent Data Quality

4.3.4.1 Availability

Availability ensures that AI data is accessible whenever needed. Reliable systems minimize downtime and ensure continuous operation, crucial for real-time AI applications.

- Evaluation Points
 - Uptime percentage meets service-level agreements (SLAs)
 - Redundancy configurations are utilized to prevent single points of failure
 - Continuous monitoring and alerting are implemented to detect access issues
- Inappropriate Examples
 - Frequent server outages disrupting data access
 - No backups, leading to inaccessible data during hardware failure
 - Delays in resolving access issues during critical tasks

4.3.4.2 Portability

Portability refers to the ability to transfer AI data seamlessly across platforms, systems, or environments without compatibility issues. It ensures flexibility and adaptability.

- Evaluation Points
 - Data is stored in widely accepted formats (e.g., CSV, JSON)
 - Standardized APIs are used for data exchange
 - Adequate documentation for data migration is in place
- Inappropriate Examples
 - Proprietary formats requiring specialized software
 - Inconsistent data structures across systems
 - Lack of metadata, causing misinterpretation during migration

4.3.4.3 Recoverability

Recoverability focuses on the system's ability to restore AI data quickly and accurately after unexpected disruptions or failures, ensuring minimal data loss.

- Evaluation Points
 - Backups are performed on a regular basis with multiple restore points maintained
 - Disaster recovery plans are tested periodically
 - The storage system must be configured with redundancy to prevent single points of failure
- Inappropriate Examples
 - Outdated backups causing irreversible data loss
 - Recovery processes requiring significant manual intervention
 - Failure to test recovery plans, leading to delays

4.3.5 Additional Data Quality Characteristics for AI/ML

4.3.5.1 Auditability

Auditability ensures data traceability and the ability to review the processes of data collection and usage. It enables accountability and compliance with ethical and legal standards. Unlike Traceability, Auditability focuses on whether data handling can be reviewed and verified.

- Evaluation Points
 - Data sources are clearly documented
 - Data collection processes are well-defined
 - Data lineage records are available
- Inappropriate Examples
 - Missing metadata for data sources
 - Ambiguous data provenance

- Inaccessible logs for key data processes

4.3.5.2 Balance

Balance ensures that data represents all relevant categories or outcomes proportionally, minimizing bias in AI models.

- Evaluation Points
 - Representation across categories is appropriately balanced for the intended use
 - Over- and under-sampling are avoided
 - Consistency is maintained across datasets
- Inappropriate Examples
 - Gender imbalance in a dataset
 - Skewed representation of geographical regions
 - Over-representation of one age group

4.3.5.3 Diversity

Diversity ensures that datasets include a wide range of perspectives, scenarios, and variations for better generalization in AI systems.

- Evaluation Points
 - Different cultural contexts are covered
 - Various scenarios and demographics are included
 - A wide range of linguistic expressions is included
- Inappropriate Examples
 - Excluding minority dialects
 - Homogeneous data in multilingual settings
 - Ignoring varied environmental factors

4.3.5.4 Effectiveness

Effectiveness refers to whether the dataset is usable for a specific AI task. It focuses on whether the data meets task requirements.

- Evaluation Points
 - Input is of sufficient quality, such as image resolution
 - There are enough samples per category
 - Usable labels and annotations are provided
- Inappropriate Examples
 - Too few major samples
 - Missing task-critical inputs

- Excessive input noise

4.3.5.5 Identifiability

Identifiability refers to whether individuals can be identified from the data. This can happen directly or through combinations of attributes. Unlike Confidentiality, Identifiability focuses on whether individuals can be identified from the data.

- Evaluation Points
 - No direct identifiers are included
 - Linkage risk is low
 - Re-identification risks are assessed
- Inappropriate Examples
 - Retention of identifiable personal details
 - Incomplete data masking
 - Reversible anonymization techniques

4.3.5.6 Relevance

Relevance refers to whether the data fits the intended use. It should include meaningful information and exclude unnecessary information. Unlike Effectiveness, Relevance focuses on whether the data is meaningful for the intended use.

- Evaluation Points
 - Relevant features are included for the AI model's domain
 - Irrelevant variables are excluded
 - Redundant information is avoided
- Inappropriate Examples
 - Off-scope records mixed in
 - Including unrelated features
 - Excessive focus on non-critical variables

4.3.5.7 Representativeness

Representativeness refers to whether the data reflects the real-world conditions or scenarios the AI model will encounter. A large gap from production conditions can reduce model performance. Representativeness can differ from Balance, because real-world data is not always evenly distributed.

- Evaluation Points
 - The dataset aligns with the characteristics of the target population
 - Expected conditions are covered
 - Sampling bias is avoided

- Inappropriate Examples
 - Over-sampling urban populations in national studies
 - Ignoring rare but critical conditions
 - Incomplete geographic coverage

4.3.5.8 Similarity

Similarity refers to whether too many samples are overly similar. Excessive similarity can reduce generalization.

- Evaluation Points
 - Sample variation is adequate
 - Near-duplicates are limited
 - No dense pattern concentration exists
- Inappropriate Examples
 - Too many duplicate samples
 - All measurements from the same equipment
 - Mostly rehashed versions of similar documents

4.3.5.9 Timeliness

Timeliness refers to whether data becomes available in time for use. Even correct data can lose value if it arrives too late. Unlike Currentness, Timeliness focuses on whether data becomes available in time, not whether it is recent enough.

- Evaluation Points
 - Data arrives on time
 - Arrival delay is within an acceptable range
 - The update frequency is suitable
- Inappropriate Examples
 - Data arriving after the deadline for necessary decisions
 - Wide variations in delay
 - The too slow update cycle

4.3.6 Sensor Data Quality Characteristics

4.3.6.1 Sensor Data

Sensor data refers to data obtained in digital form by measuring physical quantities such as temperature, humidity, location, and acceleration using sensors. It is mainly collected and accumulated in real time and utilized for device control, analysis, and service provision, forming an essential information infrastructure in modern society where IoT and smart systems are increasingly widespread.

Although each piece of sensor data is small, when aggregated in real time, it becomes a large volume of data. Many of these data are embedded in devices and services and are directly related to safety. In addition, sensors are used in diverse environments such as outdoors and inside equipment, and differences in installation conditions (e.g., measurement height, device used, measurement method) and external environmental influences can cause variability in the data. Therefore, data collection management according to sensor type and usage conditions is essential. On the other hand, when many sensors are deployed, redundancy can be utilized to compensate for individual failures using surrounding sensors or time-series data.

The quality characteristics of sensor data can be interpreted and applied based on those defined in ISO/IEC 25012 and ISO/IEC 5259-2. In addition to the format and semantic information of the collected data, it is also necessary to understand the quality derived from the state of the sensor devices that generate the data.

Based on the above, among the quality characteristics described in the previous sections, the following four characteristics are particularly important for quality management:

- Accuracy
- Completeness
- Consistency
- Precision

4.3.6.2 Device-dependent Quality Metrics

As described above, the quality of sensor data is significantly affected by the condition of the devices. Therefore, it is effective to define device-focused quality metrics for quality evaluation. The “Guidelines for Evaluating the Quality Level of Sensing Data” define these as follows:

Category	Device-dependent quality metric	Description
Design information	Device information	Level of understanding of the measurement principles, processing methods, etc. for the physical quantities (light, sound, etc.) input to the device
	Fault-tolerance	Level of device operation
	Durability	Level of decline in

Category	Device-dependent quality metric	Description
		serviceable parts
	Security measures	Level of implementation of security measures
	Communication stability	Level of operation without communication interruption or delay
Installation and adjustment	Appropriateness of installation method	Level of implementation of appropriate installation according to conditions
Operation and maintenance	System stability	Level of operation stability
	System environment monitoring	Level of installation status monitoring
	Appropriateness of updates	Level of appropriate software version operation

Source: Data Society Alliance (DSA), 2024, White Paper [“Study for Formulating Guidelines for Evaluating the Quality Level of Sensing Data”](#)

4.3.6.3 Anomalies Affecting Sensor Data Quality

Sensor data changes over time, and its characteristics may also vary, requiring corrective measures. ISO 8000-210 classifies types of anomalies that are useful for data analysts in evaluating data quality, detecting anomalies, and performing corrections, into those related to individual sensors and those related to multiple sensors.

(1) Anomalies in Individual Sensors

Typical anomalies that appear in the data patterns captured by a single sensor:

- Offset
 - Constant deviation from the true value.
- Drift
 - Gradual change over time.
- Trim
 - Adjustment applied to correct measurement error.
- Spike

- Sudden, short-lived jump.
- Noise
 - Random variations in the data.
- Data loss
 - Missing data points or gaps.
- Lack of data
 - Insufficient amount of collected data.
- Shift
 - Sudden change in the baseline.
- Drop or rise
 - Abrupt decrease or increase.
- Stuck
 - Repeated output of the same value.
- Bounded oscillation
 - Regular and limited fluctuations.
- Inconsistent frequency
 - Irregular data intervals (sampling rate irregularities).
- Different resolution
 - Data granularity deviates from expected specifications.
- Incorrect timestamp
 - Recorded time does not match the actual event time.
- Latency
 - Delay between occurrence of an event and its recording.

(2) Anomalies in Multiple Sensors

Anomalies that appear in the relationships among multiple related sensors:

- Dissimilarity
 - Significant differences in values among sensors that are supposed to measure the same target.
- Rule violation
 - Violation of predefined relationships or constraints among sensors (e.g., a low-temperature sensor reports a higher value than a high-temperature sensor).
- Inconsistent timestamp
 - Time information is not aligned among sensors that are expected to be synchronized.

4.3.6.4 Background Factors and Operational Management Items for Anomalies

The following are not defined as anomaly types in ISO 8000-210, but are factors that significantly affect data quality in sensor operation and multi-sensor integration. Although some of these conceptually correspond to the anomalies described above, they are treated here not as anomalies themselves but as their underlying causes and operational management targets.

1. Spatial and placement management
 - Improper sensor positioning: Inappropriate installation positions lead to situations where the measurement target cannot be accurately represented.
2. Noise and interference
 - Environmental noise: Unwanted signals from the surrounding environment degrade data quality.
 - Cross-talk: Interference between sensors leads to erroneous readings.
3. Resolution and sampling-related issues
 - Different granularity: Differences in measurement precision or resolution across sensors lead to inconsistencies.
 - Sampling rate differences: Inconsistent sampling frequencies across sensors result in incomplete or redundant data.
4. Faults and health management
 - Sensor drift: Measurement values change due to aging or degradation.
 - Stuck or inactive sensors: Sensors repeatedly output the same value or stop measuring.
 - Data loss: Missing data occurs due to failures or communication errors.
5. Calibration management
 - Bias error: Calibration issues cause offsets or scaling errors.
 - Lack of comparability among sensors results in inconsistent measurements.
6. Redundancy and data integration
 - Duplicate data: Overlapping sensor coverage results in redundant data.
 - Conflicting data: Contradictory values are produced across sensors.
7. Integration challenges
 - Different protocols: Differences in communication protocols make integration difficult.
 - Heterogeneous data formats: Differences in data formats require standardization.
8. Environmental impact

- Temperature, humidity, pressure: Differences in environmental conditions cause variations in sensor performance.

4.3.6.5 Processing with Cloud-Edge-IoT

Data collected by IoT devices may be processed at the edge (e.g., using edge AI), at data aggregation points, or in the cloud for large-scale processing. Data quality measures are therefore required at each location.

- Cloud: Detects unusual or biased data in the course of processing large volumes of data.
- Aggregation point: Data is aggregated in areas. Data conversion and integration, if necessary; some instructions may be sent to edge.
- Edge: At the edge, data processing such as data cleansing, recognition and anonymization processes are carried out. Corrections such as sensor-specific offsets may also be performed.



Figure 32. Data Processing Across Cloud, Edge, and IoT

5 Closing Remarks

5.1 Message

People are focusing on the cutting-edge technology of AI, but to ensure that we can use AI in society in a sustainable and secure manner, it is important to properly manage the quality of data.

Keep the principle of “Garbage in, Garbage out” in mind at all times, and maximize the value of AI.

6 Document Information

6.1 Publishing Organizations and Contributors

Japan AI Safety Institute (J-AISI) is the government initiative responsible for ensuring safety as a foundation for accelerating AI and AI-based innovation. Information-

technology Promotion Agency (IPA) is a government-funded organization focused on digital technologies and participates in J-AISI's activities.

This document has been jointly produced by J-AISI's standards team and a team of data experts from IPA's Digital Infrastructure Center, with valuable contributions from J-AISI's Data Quality Sub-Working Group.

This is a living document, and we welcome feedback from readers.

6.2 References

- ISO/IEC 25012: SQuaRE, Data quality model
- ISO/IEC 25024: SQuaRE, Measurement of data quality
- ISO 8000: Data quality
- ISO/IEC 5259: Data quality for analytics and machine learning (ML)
- ISO/IEC 8183: Data life cycle framework
- ISO/IEC 38505-1: Information technology — Governance of IT — Governance of data
- ISO 19157: Geographic information, Data quality
- DAMA-DMBOK (2nd edition, 2017), DAMA International
- Data quality guidebook for data sharing, Cabinet Office, Japan
 - https://www.chisou.go.jp/tiiki/kokusentoc/supercity/supercity_230926_guidebook.html
- White Paper “Study for Formulating Guidelines for Evaluating the Quality Level of Sensing Data”, Data Society Alliance (DSA), Japan
 - <https://data-society-alliance.org/survey-research/data-quality-evaluation-standards/>
- Machine Learning Quality Management Guideline 4th Edition, National Institute of Advanced Industrial Science and Technology (AIST), Japan
 - <https://www.digiarc.aist.go.jp/publication/aiqm/>

6.3 Version History

- Version 1.02 (2026-05-14)
 - Converted the document format from a presentation-based format to a text-based format to improve accessibility, usability, and maintainability.
 - Reorganized the document structure, including headings and hierarchy, to ensure consistency and clarity.

- Adjusted and added explanatory text to enhance readability in the text-based format. Also adjusted lists, tables, figures, and other elements.
- Made minor revisions to maintain consistency and improve the accuracy of descriptions, without changing the overall content or intent.
- Improved the accuracy of descriptions related to data quality characteristics.
- In addition to the updates described above, a machine-translated Japanese reference version has been published. Manual corrections were kept to a minimum, addressing only essential semantic inconsistencies.
- Version 1.01 (2025-12-02)
 - Revised wording and adjusted page order to improve accuracy and readability without affecting the overall content.
 - Added “About This Guide” before the main sections.
 - Content previously listed as “IV. Implementation” in Version 1.00 contained only limited material, so the relevant parts were integrated into Section III.
 - In future updates, the Implementation section will be reorganized and expanded.
- Version 1.00 (2025-03-31)
 - Released initial version.

6.4 About the Japanese Reference Translation

This English version of the guidebook is the official edition. The Japanese version is provided as a reference translation to support understanding of the content.

7 Appendix 1

7.1 Machine Learning Quality Management Guideline 4th Edition

“The Machine Learning Quality Management Guideline” published by AIST defines quality criteria and target levels for machine learning systems, with the aim of improving AI quality and reducing risks. It is primarily intended for providers of AI-enabled products and services and system developers, and aims to promote appropriate transactions and the evaluation of high-quality systems through the sharing and visualization of quality-related understanding.

This guideline organizes quality into three layers: “quality in use,” “external quality,” and “internal quality,” and positions the improvement of internal quality as leading to the achievement of external quality, which in turn enables the realization of quality in use. The following two figures illustrate the structure for achieving quality in this guideline, as well as the structure and characteristics of internal quality. Data quality is addressed as part of internal quality.

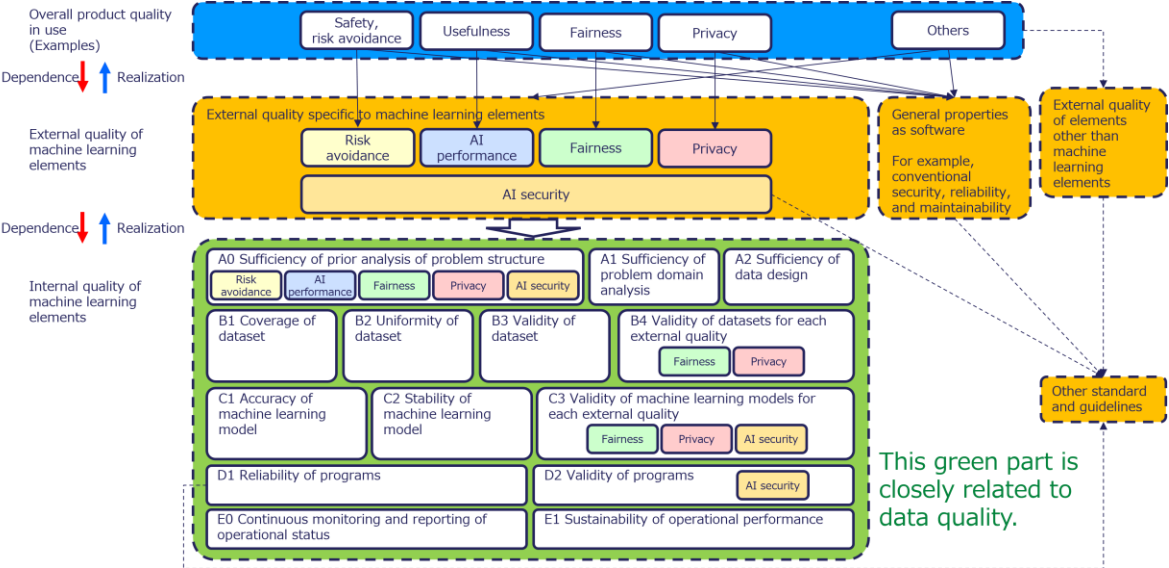


Figure 33. Structure for Achieving Product Quality

Data is a core element in AI. This guideline offers guidance on data quality management issues and responses to consider when promoting machine learning.

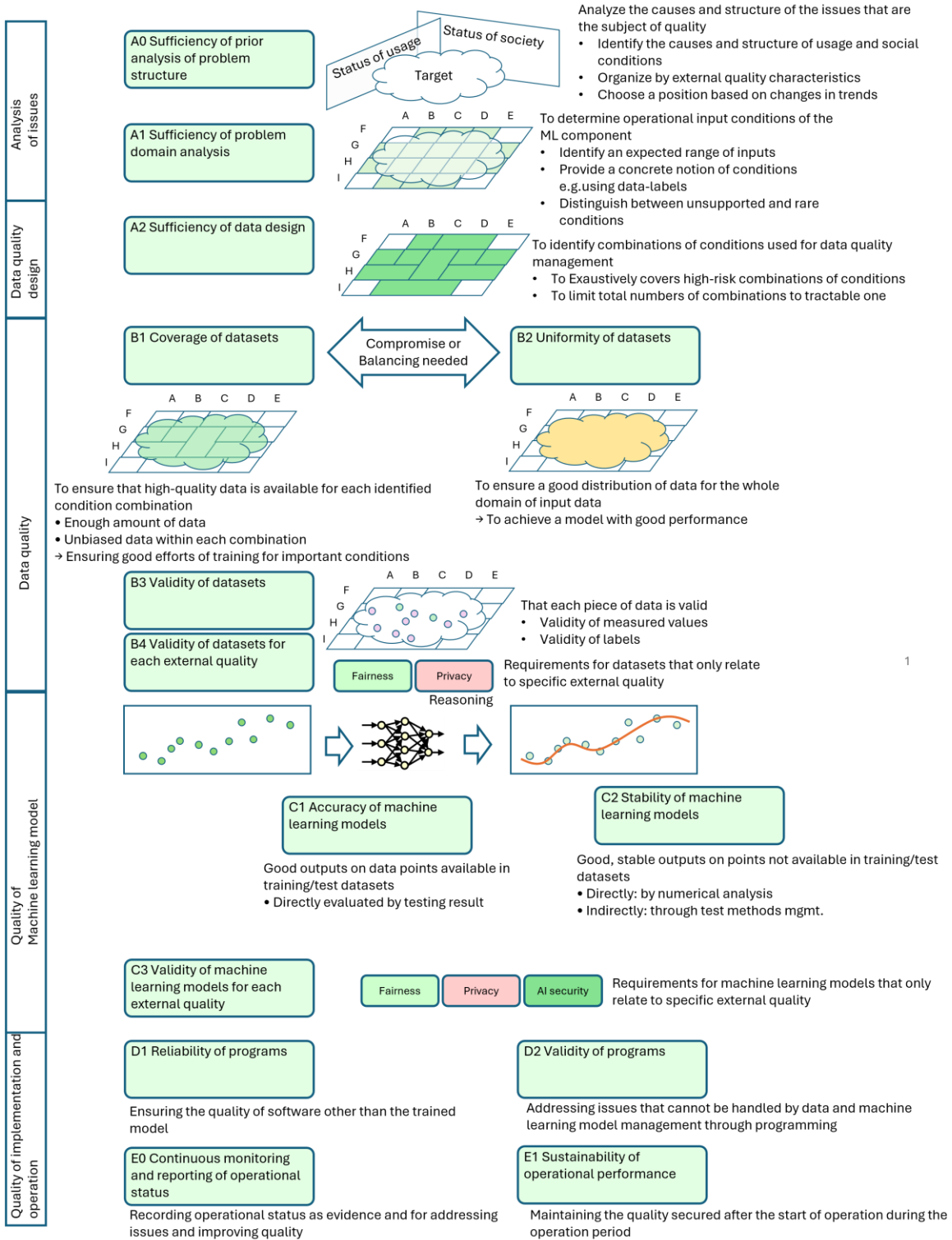


Figure 34. Internal Quality Characteristics

Source: National Institute of Advanced Industrial Science and Technology (AIST), 2023, [“Machine Learning Quality Management Guideline”](#)

8 Disclaimer

While every effort has been made to ensure the accuracy and reliability of the information provided in this guidebook, we make no warranties, either express or implied, regarding its contents. Neither organization shall be liable for any loss, damage, or claim arising from the use of this guidebook. All information is provided “as is,” and the readers assume full responsibility for any actions taken based on its content.