**Computer & Communications Industry Association**
Open Markets. Open Systems. Open Networks.
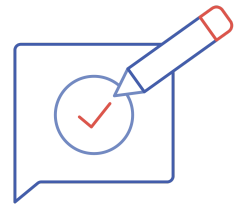
ccianet.org • @CCIAeurope

Europe

**UNDERSTANDING ARTIFICIAL INTELLIGENCE**
# Data Governance & AI – Explained

*This explainer on data governance is part of CCIA Europe's 'Understanding AI' series, which aims to inform EU policymakers and the wider public about important concepts related to artificial intelligence (AI) and the EU regulatory framework.*

## What is data governance for AI?

Data governance helps AI systems work effectively and fairly by making sure the data they use is reliable and safe. This involves having rules and processes in place to **ensure proper management of the data that is used to train and operate AI systems**. For example, it might include oversight of how data is collected, stored, accessed, and used with a view to making sure it is **secure and does not lead to biased or discriminatory outcomes**.

## Why is data governance so important?

Data governance ensures that an AI model's decisions and predictions are transparent, understandable, and avoid biases or discrimination. It involves monitoring the performance in real-world applications to detect any adverse effects and correct them if necessary.

## How does data governance work?

AI data governance involves specific steps and measures to manage data throughout the entire AI process, including development and deployment. These steps include:

1. **Data collection**: deciding what data is needed for the AI system to learn and make decisions. This involves defining the types of data to collect, where to get it from, and how to ensure data quality.

2. **Data storage and security**: storing data securely so that it can be accessed by AI systems when needed, but protecting it against unauthorised access or breaches.

3. **Data preprocessing**: cleaning and preparing the data before feeding it into the AI system. This must be done carefully to remove errors or biases that could affect the AI's performance.

4. **Model training**: training the AI model from the right data, using a well-documented training process.

5. **Model deployment**: ensuring that the data the AI model interacts with is suitable, safe, and compliant with existing regulations.

6. **Ongoing monitoring and maintenance**: continuously monitoring the AI system's performance and data usage to address potential issues or biases that may arise.

It is **crucial to make a clear distinction between input data and output data**. Input data refers to the information fed into an AI system as the initial source of knowledge, while output data is the result or response generated by the AI system based on the input data.

This distinction is **particularly relevant for data governance and the sharing of responsibilities along the AI value chain**. For input data, good governance means ensuring that the input data used to train AI models is of high quality, relevant, and representative of the real-world scenarios the AI system will encounter. For output data, governance focuses on the reliability, accuracy, and fairness of the AI system's outputs.

## How to prioritise good data governance in AI policy?

**Data governance is an important component of well-designed AI policies**. In order to avoid potential biases and discriminatory outcomes, AI systems must be trained on high-quality, representative data.

The EU's **AI Act proposal rightly includes specific data and data governance requirements for high-risk AI systems and for so-called foundation models**.

However, there is a risk that **applying the proposed stringent data governance requirements too broadly would fail to account for the differences** between the systems, the different contexts they are used in, and the variety of actors deploying, developing and using them. Currently, the burden of complying with the AI Act's data governance requirements would largely lie on developers rather than users.

However, while developers of AI systems often manage the data on which the system is initially trained, many systems ingest data from users as part of their operators. This **makes the user, rather than the provider, best placed to apply data governance requirements** as they often control how the data is retained, used and deleted.

Data governance **requirements must be technically feasible, and take into account the context of deployment** of the system. In many instances, governance requirements should apply to the deployers or users of AI systems, and not to the developers. The latter often simply don't know in what context an AI system will be used when they develop it, especially when it comes to general purpose AI (GPAI) systems that can be used for an almost unlimited number of applications.