# Digital and data: Terminology and relevant bodies

## TERMINOLOGY

**Algorithm**
An algorithm is a set of step-by-step instructions. IN Artificial Intelligence (AI), the algorithm tells the machine how to find answers to questions or solutions to a problem.

**Algorithmic bias**
Algorithmic bias can arise from the people who develop the algorithm introducing inherent prejudice in the data used to develop or train it, and/or into the application of any results. Without input from diverse teams and rigorous testing, it can be too easy for people to introduce subtle, unconscious biases, which are then perpetuated. This can mean the algorithm inappropriately advantages one group over another: e.g. it can reinforce systemic prejudices around race, ethnicity, gender, sexuality or disability.

**Algorithmic tool**
An algorithmic tool is a product, app or device that supports or solves a specific problem using complex algorithms.

**Algorithmic Transparency Standard** (ATS)
The ATS, part of the government's National Data Strategy, helps public sector organisations provide clear information about how they arrived at algorithm-assisted decisions.

**Algorithmic transparency** means being open about how algorithmic tools support decisions, providing information in a comprehensive, open, understandable, accessible and free format.

**Application Programming Interface (API)**
An API is a way for two or more computers to communicate with each other: a form of software providing a service to other pieces of software. An API opens up one organisation's applications' data to external third parties.

**Artificial Intelligence (AI)**
In its simplest form, AI combines computer science and data sets to enable problem solving: machines are taught to mimic human decision-making and learning behaviours. AI also covers fields such as *machine learning* and its sub-field *deep learning* (see below). What is called 'Weak AI' (or 'Artificial Narrow Intelligence') is AI trained to perform specific tasks (e.g. Amazon's Alexa): this drives most of the AI that currently exists. "Strong AI' is a theoretical form of AI where a machine would have self-consciousness and an intelligence equal to humans', while 'Artificial Super Intelligence' (or super-intelligence') would surpass the intelligence and ability of the human brain.

In health care, AI is mostly used to search and analyse vast data sets e.g. from health records, clinical trials, population studies and insurance claims in order to discover patterns and insights. It's claimed that AI is particularly suited to meet healthcare challenges (especially medical imaging) as it thrives on an overabundance of data, 'sees' what most humans would miss, and integrates and helps enhance existing systems.

- o *Machine learning*
With machine learning, a computer system gains the ability to use its neural network to process incoming data, identify patterns and make decisions with minimal human direction.

o *Deep learning*

Deep learning enables a computer system to process massive amounts of data quickly and prioritise the critical criteria for reaching a decision. This is a type of machine learning in which systems can accomplish complex tasks by using multiple layers of choices based on the output of the previous layer, creating increasingly 'smarter' and more abstract conclusions. Examples include facial recognition or recognising an anomaly in a medical scan.  AI can also provide healthcare organisations with algorithms to inform clinical and planning decisions, or improve service experience. Compared with machine learning, deep learning requires less human intervention.

See here for more explanation…

**Blockchain**

Blockchain is defined as 'a distributed, decentralised data ledger', but put more simply is a shared database. It was first developed through Bitcoin to disrupt the banking sector and its use of centralised ledgers. The technology enables the creation of digital records and their secure sharing and management over a network. As the name suggests, data is organised in packages or 'blocks', which are then chained together. The blocks in the chain cannot be edited: the technology only allows more blocks to be added. Each member of the network has an identical copy of the database.

Blockchain has the potential to revolutionise the handling and sharing of medical records. Its ability to record digital events, enable peer-to-peer sharing between parties, and exchange complex health data/information between patients, providers, payers and other sources (employee wellness programmes, wearable health monitors, etc.) could overcome interoperability challenges in current health IT systems. Blockchain could have an important role in claims adjudication and billing management and reduce administrative costs – e.g. for a private insurance-based system. It is also said to improve integrity in the reporting of clinical trials and population health research by reducing bias. (Studies show that roughly 30%-50% of clinical trials currently go unpublished, and that those trials that are published tend to be those with favourable results.)

While blockchain might address some challenges facing healthcare system, there are many caveats in its use concerning, for example, privacy, data-protection, and governance of blockchain networks. Plus, as the volume of data increases, new mechanisms of data storage such as data-lakes will be needed.

**Cloud technology**

'Cloud' refers to a virtual space or online platform where users can store large files and applications on remote servers, instead of data being saved directly on a user's own computer hard drive.

**Control of Patient Information (COPI)**

Health Service Control of Patient Information (COPI) Regulations 2002 allow the processing of confidential patient information for specific purposes, e.g. in relation to threats to public health. 'Processing" includes the use, dissemination and obtaining of information, the recording and holding of information, the retrieval and combination of information and the blocking and destruction of information.

In March 2020 the Government brought in emergency measures under COPI that required NHS Digital to share confidential patient information with organisations entitled to process this for Covid-19 purposes. The measures were due to expire in June 2022. However *Data Saving Lives* (2022) suggests that (COPI 2002) will be amended to facilitate access to health and care data beyond any emergency.

**Dataset**
A dataset is a collection of numbers or words that can be analysed to gain information. Datasets are often generated and stored in a tabular format according to different variables (such as height or age) with each row corresponding to a different entry (e.g. a different person). The data may come from real-life observations or measurements, or can be generated artificially (see 'synthetic data').


**Data Institutions**
Data institutions are organisations whose purpose involves stewarding data on behalf of others, often towards public, educational or charitable aims. They already play a number of vital roles, including:

- **holding data on behalf of an organisation, person, or group**, and sharing it with others who want to use it for a particular purpose.
- **combining or linking data from different sources**, and providing insights and other services back to those that have contributed data.
- **creating open datasets that anyone can access, use and share** to further a particular mission or cause.
- **developing and maintaining a common data infrastructure** for a sector or field, such as by registering identifiers or publishing open standards.

There are also a [variety of data institutions](#) emerging to support people and communities to take a more active role in stewarding data about themselves. These include data co-ops, data unions, data coalitions, and bottom-up data trusts.

**Data Processing Agreements (DPAs)**

DPAs are required under UK GDPR when a Data Controller outsources data processing to a third party vendor or partner (with 'processing' meaning anything you can possibly do with someone's personal information, such as collecting, storing, monetising or destroying it). A DPA is a legal contract between the data controller and data processer, guaranteeing that the data processor will handle the data provided appropriately in accordance with GDPR rules.

**Data sharing agreements (DSAs)**
Data sharing agreements are not mandatory but the [Information Commissioners Office](#) (ICO) considers it good practice to have a DSA in place when data is passed from one organisation (e.g. government departments or other public bodies) to another. DSAs set out the purpose of data sharing, cover what happens to the data at each stage, set standards, set out the legal basis for data sharing and help all the parties involved in sharing to be clear about their roles and responsibilities. Having a DSA in place helps parties to demonstrate they are meeting their accountability obligations under the UK GDPR.

DSAs vary depending on the scale and complexity of data sharing, for example, whether they cover the sharing of personal data or commercially sensitive information. They can take the form of an overarching document such as a Memorandum of Understanding, be one or more of a tier of documents, or part of a contract.

**Data Trusts**
[A data trust](#) is a form of data institution authorised to look after and make decisions about data, for the benefit of a wider group of stakeholders, weighing up, for example, whether

applications to access data banks pose a significant risk to an individual's privacy, and if access will strike a balance between public good, scientific discovery and value generation. With data trusts, the independent person, group or entity stewarding the data takes on a fiduciary duty. In law, a fiduciary duty is considered the highest level of obligation that one party can owe to another – in this context it involves stewarding data with impartiality, prudence, transparency and undivided loyalty.

Although data trusts are a fairly new concept and a global community-of-practice is still growing around them, there are some existing examples of independent, fiduciary stewardship of data. For example, UK Biobank was set up in 2006 to steward genetic data and samples from 0.5m people. It takes the form of a charitable company with trustees. According to the Open Data Institute, there shouldn't be inherent trust in a data trust; much depends on its trustees. There are particular concerns about data trusts and the global South: The International Digital Health and AI Research Collaborative, for example, fears that data may be handed to western researchers with no clear route to ensure that those who generated the data will benefit.

### Digital twin

A digital twin is a computer model that's an exact copy of an object in the real world (such as a bridge, a city or a biological system). Analysing the model's output can tell researchers how the real object will behave, so helping to improve its real-world design. Digital twins are key building blocks of the metaverse.

In medicine, models can be made of body systems to research the effects of various interventions (see "machine learning").

### Disruptive technologies

A disruptive technology is an innovation, possibly unproven, that significantly alters the way of doing things, or may even sweep away the systems it replaces. Disruptive technologies generally originate in start-ups and young, risk-taking companies: they contrast with 'sustaining technologies' that depend on incremental improvements in technology that already exist. Recent examples of disruptive technologies include Artificial Intelligence (AI), virtual or augmented reality, and the Internet of Things.

### Federated data platforms (FDPs)

According to NHS England,

> "Data federation is a software process that allows multiple databases to function as one. The virtual database takes data from a range of sources and converts them to a common model, providing a single source of data for front-end applications. This can facilitate access to sensitive health data, offering a potential solution to address the issue of siloed health data and barriers to data sharing."

It's stressed that an NHS FDP will focus on *connecting* datasets, rather than *collecting* data, in different 'instances' that only relevant staff can access: for example, linking a Trust's data on theatre use, waiting lists and staff rotas could help to prevent last minute cancellations of surgery.

In April 2022 NHSE announced procurement for the development of an NHS FDP (estimated value £360 million), comprised of
- the FDP itself, with ICS integration and consultancy and communications support for ICS implementation and adoption; and later
- privacy-enhancing technology
- 'a market place for applications and training
- deployment support.

The FDP will be built on 5 use cases:
- population health and person insight

- care coordination (ICS)
- elective recovery (Trust)
- vaccines and immunisation
- supply chain.

The front runner for the main FDP contract is widely assumed to be US software firm Palantir Technologies, which partnered with NHSE, Microsoft and Amazon to develop a data platform for Covid-19 response (£23.5 million for Palantir)  (challenged by Foxglove).

It is unclear how such a FDP fits with proposals to create TREs as the means of improving data privacy. Some of the FDP's data will contain identifiable information about patients.

**Fourth Industrial Revolution (4IR)**
The term 4IR has been in use from the 2000s onwards to describe the era marked by new capabilities dependent on existing technologies, such as digital systems, but in entirely new ways where technology becomes embedded in societies and even human bodies. Examples include genome editing, new forms of machine intelligence and approaches to governance relying on methods to prevent access to private messages, such as blockchain.

**Interoperability**

Interoperability is the ability of different information systems, devices and applications to access, exchange, integrate and collaboratively use data in a coordinated way, within and across organisational, regional and national boundaries and without going through a lengthy and costly consolidation process. In the context of the private sector, interoperability maximises the business value of data. In the NHS, interoperability is essential to NHS Integrated Care Systems although, currently, databases and software are not standardised.

**Knowledge economy**
The knowledge economy refers to the changing economic basis of the developed world as a result of the rapid expansion of knowledge and the increasing reliance on computerization, big data analytics, and automation. It's claimed that such an economy is dependent on intellectual capital and skills, and less dependent on traditional production processes. However, this overlooks the dependence of the knowledge economy on physical processes such as precious metal mining, semiconductor production, and computer manufacturing. In addition, the knowledge economy has created the conditions for a new form of piecework, typified by Amazon Mechanical Turk, that draws on a global, on-demand, 24/7, out-sourced workforce to carry out on-line work. This often takes the form of low-skilled, low-paid, monotonous tasks that computers are currently unable to do.

**Machine learning (ML)**
Machine learning is a field of AI using computer algorithms that can 'learn' by finding patterns in sample data. The algorithms then apply these findings to new data to make predictions or, for example, translate text or guide a robot in a new setting. In the context of healthcare, machine learning algorithms can identify tumours in scans for example.
Combining ML with a variety of data sets and almost unlimited computing power, clinical researchers can reconstruct the underlying mechanisms of disease to predict the outcome for particular difficult clinical scenarios better that usual statistics. ML is likely to play an increasing role in recruiting patients to clinical

trials by exploring patient medical records to find people who would fit trial criteria. Work in progress promises many improvements – and many difficult ethical issues - ahead.

**National Data Opt-Out**
The National data opt-out was introduced in May 2018 following recommendations from the National Data Guardian. It indicates that a patient does not want their confidential patient information to be shared for purposes beyond their individual care across the health and care system in England. It allows individuals to set a national data opt out, or reverse a previously set opt-out. It replaced the previous type 2 opt-outs that patients registered via their GP practiced. Previous type 2 opt-outs have been converted to national data opt-outs. ([https://digital.nhs.uk/data-and-information/publications/statistical/national-data-opt-out#latest-statistics](https://digital.nhs.uk/data-and-information/publications/statistical/national-data-opt-out#latest-statistics) June 2022)

**Neural network**
A neural network is a sub-set of machine learning that is at the heart of deep learning algorithms. Their name and structure are inspired by the human brain. Artificial neural networks rely on training data to learn and improve their accuracy over time. Once fine-tuned, they allow data to be classified and clustered at great speed. One of the most well-known of neural network is Google's search algorithm.

**Open data**
Open data is data that can be freely used, re-used and redistributed by anyone, subject to share alike and attributing source.  It means

- the data must be available as a whole, in a convenient and modifiable form, and at no more than a reasonable reproduction cost, preferably by downloading over the Internet.
- the data must be provided under terms that permit re-use and redistribution including the intermixing with other datasets.
- there must be universal participation: everyone must be able to use, re-use and redistribute the data without discrimination against persons or groups, or fields of endeavour. For instance, 'non-commercial' restrictions are not allowed.

Open data should not include data on specific individuals. It also relies on interoperability, which ensures different data sets can be combined.

**Metaverse**
The metaverse (sometimes confused with Web3) is potentially the next iteration of the Internet, enhances and upgraded to deliver 3D content, spatially organised information and experiences and real-time synchronous communication. Different tech companies offer different descriptions of the metaverse, but four key trends have been identified:
- 3D computing (e.g. 3D graphics, artificial visual worlds, augmented reality, digital twin);
- immersive technologies (e.g. virtual reality headsets, augmented reality glasses, spatial audio, haptics [i.e. the use of tech generated movement that the user experiences through touch));
- interconnection (e.g. interoperability of different devices, software and platforms); and
- user-generated experiences (e.g. real-time user interactions).

**Model**

A model in this context represents what has been learnt by a machine-learning algorithm – in other words, it's the output of an algorithm once it's been trained on a data set. It combines the rules, numbers and any other algorithmic-specific data structures needed to make predictions.

**Platform**

Platforms have been seen by some as the [critical hinge](#) in the history of the Internet, enabling big tech companies to unlock its profit potential. Originally 'platform' referred to something that developers build apps on top of, such as an operating system, but it now refers to different kinds of software that run on-line, especially those used by the biggest tech companies. The term suggests an aura of openness and neutrality, implying a supportive role, but actually obscures Big Tech's control of our digital spaces. The policy paper Data; Saving Lives (2022) states that modernising NHS data architecture "requires us to think of the health and care system as a platform".

**Non-fungible tokens (NFTs)**

An NFT is a token that represents the record of ownership of a unique digital asset, such as a 'jpeg' or video file. If stored and authenticated on a block chain, they are nigh impossible to tamper with because each node in the blockchain retains a complete copy of the ledger, confirming ownership. (A blockchain node refers to a stakeholder/device that plays role in governing the software of the blockchain, in place of a central entity. Their primary function is to maintain consensus of a public leger, by validating transactions and monitoring live activity to ensure a system's security.

**Privacy enhancing technologies (PETs)**

PETs are any technical methods for protecting the privacy of personal or sensitive information. This can include relatively simple technologies, such as ad-blocking browser extensions. There is particular interest in a narrower set of young, emerging PETs that are increasingly used to support secure data processing, data sharing and privacy-preserving machine learning. The most prominent PETs include

- Homomorphic encryption, which allows computations to be performed on encrypted data
- Trusted execution environments (TREs), which can protect code and data in a processing environment that is isolated from a computer's main processor and memory
- Secure multi-partner computation, in which multiple organisations collaborate to perform joint analysis on their collective data without any one organisation having to reveal their raw data to any of the others involved
- Federated analytics, an approach for applying data science techniques by moving code to the data, rather than the traditional approach of collecting data centrally
- Differentially private algorithms, which enable population-level insights about a dataset to be derived, while limiting what can be learnt about any individual in the dataset
- Synthetic data, the generation of data that is statistically consistent with a real dataset but which can replace or augment sensitive data used in data-driven applications.

**Secure Data Environments (SDEs)**

Set to become the default route for NHS and adult social organisations to provide access to their 'de-identified' data for research and analysis.  Access to data is granted to authorised researchers and analysis of data takes place within a secure on-line platform rather than data being shared and distributed.  Users interactions are recorded and monitored, and the

information they can extract has personal identifiers removed. As no data that can be linked to an individual leaves the server, and all access to data and analysis is monitored, misuse or the number of data breaches is minimised. Trusted Research Environments are a sub-set of SDEs.

**Synthetic data**

Synthetic data is generated artificially, such as by computer simulation, rather than by real-world events. It encompasses most applications of physical modelling, such as music synthesizers or flight simulators: the output approximates the real thing but is fully algorithmically generated. Synthetic data is especially useful where privacy is important, as in healthcare as the data can retain the original data's statistical properties but with any identifying information removed. It can also train algorithms where it's dangerous to use real data, e.g. teaching a self-driving car how to deal with pedestrians. Researchers may generate synthetic data to help create a baseline for any future studies and testing.

**Training AI**

In the context of AI and machine learning projects, training refers to preparing AI to properly interpret relevant data and learn from it in order to perform a task with accuracy and in the way intended. Initially, the AI is given a set of training data (a huge data set that contains certain 'training wheels', such as tags or targets, that help it interpret the data), and is asked to make decisions based on that data. Assessment at this stage allows mistakes in learning to be spotted and corrected. The next phase concerns validating assumptions about how well the AI will perform using a new set of data, as well as dealing with unexpected variables. Finally, the AI is given a dataset without the 'training wheels' to test whether it can make accurate decisions on this unstructured information and if it performs as expected.

AI training not only relies on a lot of high-quality data but a crucial preparatory step of annotation. Annotation provides the contextual guidance necessary for the AI to interpret the data properly. So far, only humans can do data annotation.

**Trusted Research Environments (TREs)**

A Trusted Research Environment (TRE) is a secure virtual environment that researchers enter in order to work on a data set remotely, rather than downloading the data onto their own local machine. Users can extract and download the results from their analyses but it's claimed that individual patients' data always stays within the secure environment (or rather, no data linked to an individual leaves the server). TREs are diverse in their aims and design, but a robust TRE is said to help analysts to work effectively with data while also preventing and detecting misuse, and providing fully open and detailed public logs of all actions on patients' records. TREs mean that gatekeepers controlling access to data can be proportionately more permissive because it's possible to manage risk more effectively than when data is disseminated.

In addition to reducing privacy risk, good TREs can also provide a more efficient and collaborative computational environment for all data users: it has been estimated that 80% of the work for data science with NHS records is spent on data preparation; this is currently delivered in a diverse, duplicated and ad hoc fashion: different teams and individuals in different organisations or settings are using different methods and tools for even basic data management work to do the same or very similar tasks, often on the same national NHS datasets, such as GP data or Hospital Episodes Statistics. If this work is done in national or broadly standardised analytic environments, then code and working methods become portable, open, discussable, reviewable and re-usable.

[A review of TREs](#) led by Ben Goldacre claims that TREs represent an opportunity to modernise the data management and analysis work that's done across a system, with benefits including:

- the replacement of hundreds of different analytic siloes, data centres and working practices with a small number of more standardised environments that facilitate more modern, efficient approaches to data science.
- a reduction in the number of data centres, and therefor a reduced number of cost centres.
- fewer 'attack surfaces' for hackers.
- the creation of working environments where a number of expert software developers can assist all colleagues, packaging up the code for recurring tasks into adequately documented "functions" and "libraries" for easy re-use.
- the collaborative development of effective interactive data tools for less skilled users.
- ensuring that all data curation code is shared with all subsequent users for review, validation, re-use, and iterative modification
- making modern, open, collaborative, computational approaches to data analysis the norm.

**Web 3.0**

For some, 'web3' has become a general term for any emerging digital concepts relating to blockchain and crypto. It is also used to refer specifically to the third (and future) version of the Internet.

Web 1.0 (or the "world wide web') (1991 – 2004) provided access to static, read-only information on web pages. Users could only passively browse texts, pictures and short videos.

Web 2.0 (2005-present), emerged in parallel with social media, provided access to dynamic and fluid media, user-generated content, interoperability, and interactivity. Underpinned more recently by the use of metadata, such as cookies and an advertising-revenue driven model.

Web 3.0 (future) has no single definition. For some, web3 no longer relies on large technology companies to run the Internet or apps, with access dependent on providing your data. Instead, the web becomes decentralised, using apps that run on blockchain and/or peer-to-peer network technology. Users become joint owners/controllers of the Internet and the data held on it, and are financially incentivised to create, govern or contribute to projects. The web is also connected with home devices, wearable devices etc. which means it connects everything – it's the Internet of Things.

Commentators associated with the Open Rights Group argue that Web 3.0 will provide increased data security, scalability, and privacy for users and combat the influence of large technology companies. But they also raise concerns about the decentralised web component of Web 3.0, citing the potential for inadequate moderation and the proliferation of harmful content.

## RELEVANT ORGANISATIONS AND BODIES

[Ada Lovelace Institute](#), established by the Nuffield Foundation, aims to ensure the benefits of data and AI are justly and equitably distributed, and must enhance individual and social wellbeing.

Alan Turing Institute is the national institute for data science and artificial intelligence, funded by grants from research councils, university partners and other partnerships (e.g. Roche, Siemens, Microsoft)

*Centre for Data Ethics and Innovation* is part of the Department for Digital, Culture, Media and Sport and described as an expert body enabling the trustworthy use of data and AI.

*Centre for Improving Data Collaboration* - a business unit within NHSX/NHS Transformation Directorate that provides specialist commercial and legal advice to NHS Trusts, medical charities and 'other' health sector organisations entering into data partnerships. The core aim is to ensure fair value return to the health system for sharing data assets, with health data owners "feeling empowered to seek out commercial relationships".

*Independent Group Advising on the Release of Data* (IGARD) has two purposes: i) to make general recommendations and observations to NHS Digital about its processes, policies and procedures, including transparency measures; and ii) to scrutinise and advise NHS Digital on the appropriateness of requests for dissemination of confidential information (including personal data).

*Data Standards Authority* establishes standards to make it easier and more effective to access and use data across government

*Information Commissioner's Office* was set up as the UK's independent body in order to uphold information rights. If the Data Protection and Digital Information Bill 2022 is passed, the ICO will become a corporate body known as the Information Commission, will loose much of its freedom from political control, and will have a new duty 'to regard economic growth, innovation and competition issues'.

*International Digital Health and AI Research Collaborative* (I-Dair) works to develop a global platform to enable inclusive and responsible research into digital health and AI for health. It aims to level the playing field and take advantage of the digital revolution to support the collaborative development of health solutions that benefit all countries and communities.

*National Data Guardian* advises and challenges the health and care system to help ensure that citizens' confidential information is safeguarded and used properly.

*Midata* – an example of a data trust organised as a co-operative, rooted in the idea that data can be used for the common good: citizens or patients (who may become members of the cooperative) have a right to decide what happens to their data, and grant selective access to their data.

*Office for Artificial Intelligence* A government department currently considering the UK's "pro-innovation approach" to governing AI with support from stakeholders.

*Open Data Institute* is a non-profit organisation that advocates for an open, trustworthy data ecosystem. It's funded by grants and commercial revenue.

*UK Health Data Research Alliance* an independent grouping of healthcare and research organisations that aims to establish best practice for the ethical use of research data. Alliance members include NHS Trusts, medical charities, universities, NHSE and partnership such as Q Research, described as a collaboration between the University of Oxford and EMIS, the leading IT provider for computer systems in primary care.