

---

---

## Information technology — Big data — Overview and vocabulary

*Technologies de l'information — Mégadonnées — Vue d'ensemble et  
vocabulaire*

**iTeh STANDARD PREVIEW**  
**(standards.iteh.ai)**

[ISO/IEC 20546:2019](https://standards.iteh.ai/catalog/standards/sist/9fe83df4-5e01-4af9-b2e4-65f5dd3410b4/iso-iec-20546-2019)

<https://standards.iteh.ai/catalog/standards/sist/9fe83df4-5e01-4af9-b2e4-65f5dd3410b4/iso-iec-20546-2019>



**iTeh STANDARD PREVIEW**  
**(standards.iteh.ai)**

[ISO/IEC 20546:2019](https://standards.iteh.ai/catalog/standards/sist/9fe83df4-5e01-4af9-b2e4-65f5dd3410b4/iso-iec-20546-2019)

<https://standards.iteh.ai/catalog/standards/sist/9fe83df4-5e01-4af9-b2e4-65f5dd3410b4/iso-iec-20546-2019>



**COPYRIGHT PROTECTED DOCUMENT**

© ISO/IEC 2019

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office  
CP 401 • Ch. de Blandonnet 8  
CH-1214 Vernier, Geneva  
Phone: +41 22 749 01 11  
Fax: +41 22 749 09 47  
Email: [copyright@iso.org](mailto:copyright@iso.org)  
Website: [www.iso.org](http://www.iso.org)

Published in Switzerland

# Contents

	Page
<b>Foreword</b> .....	<b>iv</b>
<b>Introduction</b> .....	<b>v</b>
<b>1 Scope</b> .....	<b>1</b>
<b>2 Normative references</b> .....	<b>1</b>
<b>3 Terms, definitions and abbreviated terms</b> .....	<b>1</b>
3.1 Terms and definitions.....	1
3.2 Abbreviated terms.....	6
<b>4 Key characteristics of big data</b> .....	<b>6</b>
4.1 General.....	6
4.2 Key data characteristics.....	6
4.2.1 Data volume.....	6
4.2.2 Data velocity.....	6
4.2.3 Data variety.....	6
4.2.4 Data variability.....	6
4.3 Key data processing characteristics.....	7
4.3.1 Data science.....	7
4.3.2 Data volatility.....	7
4.3.3 Data veracity.....	7
4.3.4 Benefit.....	7
4.3.5 Data visualization.....	7
4.3.6 Structured and unstructured data.....	7
4.3.7 Scaling.....	7
4.3.8 Distributed file system.....	8
4.3.9 Distributed data processing.....	8
4.3.10 Non-relational databases.....	8
<b>Annex A (informative) Cross-cutting concepts of big data</b> .....	<b>9</b>
<b>Bibliography</b> .....	<b>12</b>

## Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialized system for worldwide standardization. National bodies that are members of ISO or IEC participate in the development of International Standards through technical committees established by the respective organization to deal with particular fields of technical activity. ISO and IEC technical committees collaborate in fields of mutual interest. Other international organizations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular, the different approval criteria needed for the different types of document should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see [www.iso.org/directives](http://www.iso.org/directives)).

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO and IEC shall not be held responsible for identifying any or all such patent rights. Details of any patent rights identified during the development of the document will be in the Introduction and/or on the ISO list of patent declarations received (see [www.iso.org/patents](http://www.iso.org/patents)) or the IEC list of patent declarations received (see <http://patents.iec.ch>).

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation of the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT) see [www.iso.org/iso/foreword.html](http://www.iso.org/iso/foreword.html).

This document was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 42, *Artificial intelligence*.

Any feedback or questions on this document should be directed to the user's national standards body. A complete listing of these bodies can be found at [www.iso.org/members.htm](http://www.iso.org/members.htm)

## Introduction

The big data paradigm is a rapidly changing field with rapidly changing technologies.

The term big data implies datasets that are extensive in volume, velocity, variety and/or variability. The term does not, however, represent data that is simply larger than before, since this has happened on a regular basis for decades. The specific occurrence that has led to the widespread usage of the term big data is that in the mid-2000s, extensive datasets could no longer be handled using extant data systems. The big data techniques represented a shift at that time to use distributed data management and processing through horizontal scaling to achieve the needed performance efficiency at an affordable cost.

In the evolution of data processing systems, there have been a number of times when the need for efficient, cost-effective data analysis has forced a change in existing technologies. For example, the move to a relational model occurred when methods to reliably handle changes to structured data led in the 1980s to the shift to relational databases that modelled relational algebra. That was a fundamental shift in data handling. The revolution in technologies referred to as big data has arisen because the relational model could no longer efficiently handle all the needs for analysis of large and often unstructured datasets. It is not just that data is larger than before, as data has been steadily getting larger for decades. The big data revolution is instead a one-time fundamental shift in architecture towards parallelization, just as the shift to the relational model was a one-time shift. As relational databases evolved to greater efficiencies over decades, so too will big data technologies continue to evolve. Many of the conceptual underpinnings of big data have been around for years, but the years since the mid-2000s have seen an explosion in scaling technologies and their maturation and application to scaled data systems.

The term big data is overloaded in common usage and is used to represent a number of related concepts, in part because several distinct system dimensions are consistently interacting with each other. To understand this revolution, the interplay of the following aspects needs to be considered: the data and processing characteristics of the datasets, the analysis of the datasets, the performance of the systems that handle the data, the business considerations of cost effectiveness, and the new engineering and analysis techniques for distributed data processing using horizontal scaling.

[Annex A](#) provides an overview of several concepts from the broader computing domain which are cross-cutting with respect to big data.

**iTeh STANDARD PREVIEW**  
**(standards.iteh.ai)**

[ISO/IEC 20546:2019](#)

<https://standards.iteh.ai/catalog/standards/sist/9fe83df4-5e01-4af9-b2e4-65f5dd3410b4/iso-iec-20546-2019>

# Information technology — Big data — Overview and vocabulary

## 1 Scope

This document provides a set of terms and definitions needed to promote improved communication and understanding of this area. It provides a terminological foundation for big data-related standards.

This document provides a conceptual overview of the field of big data, its relationship to other technical areas and standards efforts, and the concepts ascribed to big data that are not new to big data.

## 2 Normative references

There are no normative references in this document.

## 3 Terms, definitions and abbreviated terms

For the purposes of this document, the following terms and definitions apply.

ISO and IEC maintain terminological databases for use in standardization at the following addresses:

— ISO Online browsing platform: available at <https://www.iso.org/obp>

— IEC Electropedia: available at <http://www.electropedia.org/>

<https://standards.iteh.ai/catalog/standards/sist/9fe83df4-5e01-4af9-b2e4-65f5dd3410b4/iso-iec-20546-2019>

### 3.1 Terms and definitions

#### 3.1.1

##### **benefit**

advantage to the organization of the actionable knowledge derived from an analytic system

Note 1 to entry: Benefit is often ascribed to big data due to the understanding that data has potential value that was typically not considered previously.

#### 3.1.2

##### **big data**

extensive *datasets* (3.1.11) — primarily in the *data* (3.1.5) characteristics of volume, variety, velocity, and/or variability — that require a scalable technology for efficient storage, manipulation, management, and analysis

Note 1 to entry: Big data is commonly used in many different ways, for example as the name of the scalable technology used to handle big data extensive datasets.

#### 3.1.3

##### **cloud computing**

paradigm for enabling network access to a scalable and elastic pool of shareable physical or virtual resources with self-service provisioning and administration on-demand

Note 1 to entry: Examples of resources include servers, operating systems, networks, software, applications, and storage equipment.

[SOURCE: ISO/IEC 17788:2014, 3.2.5]

### 3.1.4

#### cluster

<distributed data processing> set of functional units under common control

[SOURCE: ISO/IEC 2382:2015, 2120586]

### 3.1.5

#### data

reinterpretable representation of information in a formalized manner suitable for communication, interpretation, or processing

Note 1 to entry: Data can be processed by humans or by automatic means.

[SOURCE: ISO/IEC 2382:2015, 2121272]

### 3.1.6

#### data analytics

composite concept consisting of data acquisition, data collection, data validation, *data processing* (3.1.9), including data quantification, data visualization, and data interpretation

Note 1 to entry: Data analytics is used to understand objects represented by *data* (3.1.5), to make predictions for a given situation, and to recommend on steps to achieve objectives. The insights obtained from analytics are used for various purposes such as decision-making, research, sustainable development, design, planning, etc.

### 3.1.7

#### database

collection of *data* (3.1.5) organized according to a conceptual structure describing the characteristics of these *data* and the relationships among their corresponding entities, supporting one or more application areas

[SOURCE: ISO/IEC 2382:2015, 2121413]

ISO/IEC 20546:2019

<https://standards.iteh.ai/catalog/standards/sist/9fe83df4-5e01-4af9-b2e4-65f5dd3410b4/iso-iec-20546-2019>

### 3.1.8

#### data model

pattern of structuring *data* (3.1.5) in a *database* (3.1.7) according to the formal descriptions in its information system and according to the requirements of the database management system to be applied

[SOURCE: ISO/IEC 2382:2015, 2125519]

### 3.1.9

#### data processing

systematic performance of operations upon *data* (3.1.5)

Note 1 to entry: Example: Arithmetic or logic operations upon data, merging or sorting of data, or operations on text, such as editing, sorting, merging, storing, retrieving, displaying, or printing.

Note 2 to entry: The term data processing should not be used as a synonym for information processing.

[SOURCE: ISO/IEC 2382:2015, 01.01.06]

### 3.1.10

#### data science

extraction of actionable knowledge from *data* (3.1.5) through a process of discovery, or hypothesis and hypothesis testing

### 3.1.11

#### data set

#### dataset

identifiable collection of *data* (3.1.5) available for access or download in one or more formats

[SOURCE: Adapted from ISO 19115-2:2009, 4.7]



**3.1.12****data type**

datatype

defined set of *data* (3.1.5) objects of a specified data structure and a set of permissible operations, such that these *data* objects act as operands in the execution of any one of these operations

Note 1 to entry: Example: An integer type has a very simple structure, each occurrence of which, usually called value, is a representation of a member of a specified range of whole numbers and the permissible operations include the usual arithmetic operations on these integers.

Note 2 to entry: The term "type" may be used instead of "data type" when there is no ambiguity.

Note 3 to entry: Data type; datatype: terms and definition standardized by ISO/IEC [ISO/IEC 2382-15:1999].

Note 4 to entry: 15.04.01 (17.05.08) (2382).

[SOURCE: ISO/IEC 2382:2015, 2122374]

**3.1.13****data variability**

changes in transmission rate, format or structure, semantics, or quality of *datasets* (3.1.11)

**3.1.14****data variety**

range of formats, logical models, timescales, and semantics of a *dataset* (3.1.11)

Note 1 to entry: Data variety refers to irregular or heterogeneous data structures, their navigation, query, and data typing.

**3.1.15****data velocity**

rate of flow at which *data* (3.1.5) is created, transmitted, stored, analysed or visualised

**3.1.16****data veracity**

completeness and/or accuracy of *data* (3.1.5)

Note 1 to entry: Data veracity refers to descriptive data and self-inquiry about objects to support real-time decision-making.

**3.1.17****data volatility**

characteristic of *data* (3.1.5) pertaining to the rate of change of these data over time

[SOURCE: ISO/IEC 2382:2015, 17.06.06]

**3.1.18****data volume**

extent of the amount of *data* (3.1.5) relevant to impacting computation and storage resources and their management during data processing

Note 1 to entry: Data volume becomes important in dealing with large *datasets* (3.1.11), including their.

**3.1.19****distributed data processing**

*data processing* (3.1.9) in which the performance of operations is dispersed among the nodes in a computer network

[SOURCE: ISO/IEC 2382:2015, 18.01.08]

**3.1.20****distributed file system**

system which manages files and folders across multiple networked systems

3.1.21

**file**

named set of records treated as a unit

[SOURCE: ISO/IEC 2382:2015, 04.07.10]

3.1.22

**gather**

consolidation of results from multiple nodes in a cluster

Note 1 to entry: See *scatter-gather* (3.2.33).

3.1.23

**horizontal scaling**

providing a single logical unit through the connection of multiple hardware and software

Note 1 to entry: The example of horizontal scaling is increasing the performance of distributed data processing through the addition of nodes in the cluster for additional resources.

Note 2 to entry: Horizontal scaling for increasing performance is also referred to as scale-out.

3.1.24

**metadata**

*data* (3.1.5) about data or data elements, possibly including their data descriptions, and data about data ownership, access paths, access rights and *data volatility* (3.1.17)

[SOURCE: ISO/IEC 2382:2015, 17.06.05]

STANDARD PREVIEW  
(standards.iteh.ai)

3.1.25

**non-relational database**

*database* (3.1.7) that does not follow a *relational model* (3.1.31)

Note 1 to entry: NoSQL, which is typically translated as "non SQL" or "not only SQL", is the term in common usage to refer to databases that do not conform to a relational model.

3.1.26

**non-relational model**

logical *data model* (3.1.10) that does not follow a *relational model* (3.1.31) for the storage and manipulation of *data* (3.1.5)

3.1.27

**parallel**

pertaining to a process in which all events occur within the same interval of time, each one handled by a separate but similar functional unit

Note 1 to entry: Example: The parallel transmission of the bits of a computer word along the lines of an internal bus.

[SOURCE: ISO/IEC 2382:2015, 03.02.01]

3.1.28

**partially structured data**

*data* (3.1.5) that has some organization

Note 1 to entry: Partially structured data is often referred to as semi-structured data by industry.

Note 2 to entry: examples of partially structured data are records with free text fields in addition to more structured fields. Such data is frequently represented in computer interpretable/parsible formats such as XML or JSON.

**3.1.29****relational algebra**

algebra for expressing and manipulating relations

[SOURCE: ISO/IEC 2382:2015, 17.04.08]

**3.1.30****relational database**

*database* (3.1.7) in which the data are organized according to a *relational model* (3.1.31)

[SOURCE: ISO/IEC 2382:2015, 17.04.05]

**3.1.31****relational model**

*data model* (3.1.10) whose structure is based on a set of relations

[SOURCE: ISO/IEC 2382:2015, 17.04.04]

**3.1.32****scatter**

distribution of processing across multiple nodes in a *cluster* (3.1.4)

Note 1 to entry: See *scatter-gather* (3.2.33).

**3.2.33****scatter-gather**

form of processing of large *datasets* (3.1.14) where the computation required is divided up and distributed across multiple nodes in a cluster and the overall result is combined from the results from each node

Note 1 to entry: Scatter-gather processing typically requires an algorithmic change to the processing software. An example of scatter-gather data processing is MapReduce.

**3.1.34****streaming data**

*data* (3.1.5) passing across an interface from a source that is operating continuously

[SOURCE: ISO/IEC 19784-4:2011, 4.4]

**3.1.35****structured data**

*data* (3.1.5) which are organized based on a pre-defined (applicable) set of rules

Note 1 to entry: The predefined set of rules governing the basis on which the data is structured needs to be clearly stated and made known.

Note 2 to entry: A pre-defined data model is often used to govern the structuring of data.

**3.1.36****SQL**

database language specified by ISO/IEC 9075

Note 1 to entry: SQL is sometimes interpreted to stand for Structured Query Language, but that name is not used in the ISO/IEC 9075 series

**3.1.37****unstructured data**

*data* (3.1.5) which are characterized by not having any structure apart from that record or file level

Note 1 to entry: On the whole unstructured data is not composed of data elements.

EXAMPLE An example of unstructured data is free text.