

ISO/IEC JTC 1
Information technology
Secretariat: ANSI (United States)

Document type: Text for CD ballot or comment

Title: Text of CD 20546, Information Technology— Big Data— Overview and Vocabulary

Status: Please submit your vote via the online balloting system.

Date of document: 2016-04-12

Source: Project Editor

Expected action: VOTE

Action due date: 2016-07-13

Email of secretary: lrajchel@ansi.org

Committee URL: <http://isotc.iso.org/livelink/livelink/open/jtc1>



ISO/IEC JTC 1 N12986

2016-04-12

Replaces:

**ISO/IEC JTC 1
Information Technology**

Document Type: CD for ballot

Document Title: Text of CD 20546, Information Technology— Big Data— Overview and Vocabulary

Document Source: Project Editor

Project Number:

Document Status: Please submit your vote via the online balloting system.

Action ID: Vote

Due Date: **2016-07-13**

Pages:

ISO/IEC CD 20546	
Date: 2016-04-12	Reference number: ISO/IEC JTC 1 N12986
Supersedes document	

THIS DOCUMENT IS STILL UNDER STUDY AND SUBJECT TO CHANGE. IT SHOULD NOT BE USED FOR REFERENCE PURPOSES.

ISO/IEC JTC 1 INFORMATION TECHNOLOGY Secretariat: USA (ANSI)	Circulated to P- and O-members, and to technical committees and organizations in liaison for: - discussion at - comment by - voting by (P-members only) <p style="text-align: center;">2016-07-13</p> Please return all votes and comments via the online balloting system.
--	---

ISO/IEC JTC 1

Title: Text of CD 20546, Information Technology— Big Data— Overview and Vocabulary

Project: 1.20546

Introductory note:

Recipients of this document are invited to submit notification of any relevant patent rights of which they are aware and to provide supporting documentation.

1
2
3
4
5

Reference number of working document: **ISO/JTC 1/ N XXXXX**

Date: 2016-03-xx

Reference number of document: **ISO/IEC CD 20546**

Committee identification: ISO/JTC 1/WG 9

Secretariat: ANSI

6 **Information Technology— Big Data— Overview and Vocabulary**

7

8

Warning

9
10

This document is not an ISO International Standard. It is distributed for review and comment. It is subject to change without notice and may not be referred to as an International Standard.

11
12

Recipients of this draft are invited to submit, with their comments, notification of any relevant patent rights of which they are aware and to provide supporting documentation.

13

Copyright notice

This ISO document is a working draft or committee draft and is copyright-protected by ISO. While the reproduction of working drafts or committee drafts in any form for use by participants in the ISO standards development process is permitted without prior permission from ISO, neither this document nor any extract from it may be reproduced, stored or transmitted in any form for any other purpose without prior written permission from ISO.

Requests for permission to reproduce this document for the purpose of selling it should be addressed as shown below or to ISO's member body in the country of the requester:

Secretariat of ISO/IEC JTC 1/WG 9
American National Standards Institute,
25 West 43rd Street, New York, NY 10036
Tel. + 1 212 642 4904
Fax + 1 212 840 2298
E-mail jgarner@itic.org
Web <http://www.iso.org/jtc1> (public web site)

Reproduction for sales purposes may be subject to royalty payments or a licensing agreement.

Violators may be prosecuted.

31	Contents	Page
32	Foreword	iv
33	Introduction	v
34	1 Scope	1
35	2 Normative references	1
36	3 Terms and definitions	1
37	3.1 Terms Defined Elsewhere	1
38	3.2 Terms Defined in this International Standard	3
39	3.3 Abbreviations	3
40	4 Overview of Big Data	3
41	4.1 Volume	4
42	4.2 Velocity	4
43	4.3 Variety	4
44	4.4 Variability	4
45	Annex A (informative) Additional Related Concepts	5
46	A.1 Volatility	5
47	A.2 Distributed File System	5
48	A.3 Non-relational Models	5
49	A.4 Non-relational Database	5
50	A.5 Scatter-gather	5
51	A.6 Computational Portability	5
52	A.7 Data Science	5
53	A.8 Value	5
54	A.9 Veracity	6
55	A.10 Unstructured Data	6
56	Annex B (informative) Cross-Cutting Concepts of Big Data	7
57	B.1 Metadata	7
58	B.2 Analytics/Statistics	7
59	B.3 Cluster Computing	7
60	B.4 Cloud Computing	7
61	B.5 Security	7
62	B.6 Privacy	7
63	B.7 SQL	8
64	B.8 Parallel Computing	8
65	B.9 Internet of Things	8
66	B.10 Programming Languages	8
67	Bibliography	9
68		
69		

70 **Foreword**

71 ISO (the International Organization for Standardization) and IEC (the International Electrotechnical
72 Commission) form the specialized system for worldwide standardization. National bodies that are members of
73 ISO or IEC participate in the development of International Standards through technical committees
74 established by the respective organization to deal with particular fields of technical activity. ISO and IEC
75 technical committees collaborate in fields of mutual interest. Other international organizations, governmental
76 and non-governmental, in liaison with ISO and IEC, also take part in the work. In the field of information
77 technology, ISO and IEC have established a joint technical committee, ISO/IEC JTC 1.

78 International Standards are drafted in accordance with the rules given in the ISO/IEC Directives, Part 2.

79 The main task of the joint technical committee is to prepare International Standards. Draft International
80 Standards adopted by the joint technical committee are circulated to national bodies for voting. Publication as
81 an International Standard requires approval by at least 75 % of the national bodies casting a vote.

82 Attention is drawn to the possibility that some of the elements of this document may be the subject of patent
83 rights. ISO and IEC shall not be held responsible for identifying any or all such patent rights.

84 ISO/IEC 20546 was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information Technology*.

85

86 **Introduction**

87 The Big Data paradigm is a rapidly changing field with rapidly changing technologies. This standard will
88 provide the normative definitions and the vocabulary needed to promote improved communication and
89 understanding of this emerging area.

90 This international standard will provide the conceptual overview of the emerging field of Big Data, its
91 relationship to other technical areas and standards efforts, and the concepts ascribed to big data that are not
92 new to Big Data. This international standard is expected to be time sensitive and will need revision over time.

93

94 Information Technology — Big Data — Overview and 95 Vocabulary

96 1 Scope

97 This International Standard provides an overview of Big Data along with a set of terms and definitions. It
98 provides a terminological foundation for Big Data-related standards.

99 2 Normative references

100 The following International Standards contain provisions which, through reference in this text, constitute
101 provisions of this International Standard. At the time of publication, the editions indicated were valid. All
102 Standards are subject to revision, and parties to agreements based on this International Standard are
103 encouraged to investigate the possibility of applying the most recent edition of the Standards listed below.
104 Members of IEC and ISO maintain registers of currently valid International Standards.

105 2.1 Identical International Standards

106 None.

107 2.2 Paired International Standards

108 None.

109 2.3 Additional references

110 None.

111 3 Terms and definitions

112 3.1 Terms Defined Elsewhere

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

114 3.1.1

115 data processing

116 the systematic performance of operations upon data

117 NOTE 1 - Example: Arithmetic or logic operations upon data, merging or sorting of data, assembling or
118 compiling of programs, or operations on text, such as editing, sorting, merging, storing, retrieving, displaying,
119 or printing.

120 NOTE 2 - The term data processing should not be used as a synonym for information processing.

121 [SOURCE: ISO/IEC 2382-1;1999EF, 01.01.06]

- 122 **3.1.2**
123 **distributed data processing**
124 data processing in which the performance of operations is dispersed among the nodes in a computer network
- 125 [SOURCE: ISO/IEC 2382-18;1999EF, 18.01.08]
- 126 **3.1.3**
127 **data volatility**
128 characteristic of data pertaining to the rate of change of these data over time
- 129 [SOURCE: ISO/IEC 2382-17:1999, 17.06.06]
- 130 **3.1.4**
131 **cluster (in distributed data processing)**
132 A set of functional units under common control
- 133 [SOURCE: ISO/IEC 2382-18;1999EF 18.05.06]
- 134 **3.1.5**
135 **parallel**
136 Pertaining to a process in which all events occur within the same interval of time, each one handled by a
137 separate but similar functional unit.
- 138 NOTE - Example: The parallel transmission of the bits of a computer word along the lines of an internal bus.
- 139 [SOURCE: ISO/IEC 2382-3: 1987EF 03.02.01]
- 140 **3.1.6**
141 **relational algebra**
142 algebra for expressing and manipulating relations
- 143 [SOURCE: ISO/IEC 2382-17:1999, 17.04.08]
- 144 **3.1.7**
145 **relational model**
146 data model whose structure is based on a set of relations
- 147 [SOURCE: ISO/IEC 2382-17:1999]
- 148 **3.1.8**
149 **relational database**
150 A database in which the data are organized according to a relational model
- 151 [SOURCE: ISO/IEC 2382-17 17.04.05]
- 152 **3.1.9**
153 **data type**
154 defined set of data objects of a specified data structure and a set of permissible operations, such that these
155 data objects act as operands in the execution of any one of these operations
- 156 NOTE - Example: An integer type has a very simple structure, each occurrence of which, usually called value,
157 is a representation of a member of a specified range of whole numbers. The permissible operations include
158 the usual arithmetic operations on these integers.
- 159 [SOURCE: ISO/IEC 2382-17:1999, 17.05.08]
- 160 **3.1.10**
161 **streaming data**
162 data passing across an interface from a source that is operating continuously

163 [SOURCE: ISO/IEC 19784-4:2011]

164 **3.1.11**

165 **metadata**

166 data about data or data elements, possibly including their data descriptions, and data about data ownership,
167 access paths, access rights and data volatility

168 [SOURCE: ISO/IEC 2382-17:1999, 17.06.05]

169 **3.2 Terms Defined in this International Standard**

170 **3.2.1**

171 **big data**

172 extensive datasets – primarily in the characteristics of volume, variety, velocity, and/or variability – that require
173 a scalable architecture for efficient storage, manipulation, and analysis

174 NOTE - Big Data is commonly used in many different ways, for example as the name of the scalable
175 technology used to handle big data extensive datasets.

176 **3.2.2**

177 **horizontal scaling**

178 increasing the performance of distributed data processing through the addition of nodes in the cluster for
179 additional resources

180 NOTE - horizontal scaling is also referred to as scale-out

181 **3.2.3**

182 **vertical scaling**

183 increasing the performance of data processing through improvements to processors, memory, storage, or
184 connectivity

185 NOTE - vertical scaling is also referred to as scale-up

186 **3.3 Abbreviations**

187 For the purposes of this International Standard, the following abbreviations apply:

188 None.

189 **4 Overview of Big Data**

190 The term¹⁾ *big data* implies datasets that are extensive in volume, velocity, or variety. The term does not,
191 however, represent data that is simply bigger than before, since this has happened on a regular basis for
192 decades. The specific occurrence that has led to the widespread usage of the term big data is that in the mid
193 2000's, extensive datasets could no longer be handled using extant data system architectures. The new big
194 data techniques represented a shift at that time to use *distributed data processing* through *horizontal scaling*
195 to achieve the needed performance efficiency at an affordable cost.

196 In the evolution of data systems, there have been a number of times when the need for efficient, cost effective
197 data analysis has forced a change in existing technologies. For example, the move to a *relational model*
198 occurred when methods to reliably handle changes to structured data led in the 1980's to the shift to *relational*
199 *databases* that modelled *relational algebra*. That was a fundamental shift in data handling. The revolution in
200 technologies referred to as *big data* has arisen because the *relational model* could no longer efficiently handle
201 all the needs for analysis of large and often unstructured datasets. It is not just that data is bigger than before,
202 as data has been steadily getting larger for decades. The big data revolution is instead a one-time

1) Note that all terms referenced in Clause 3 are italicized.

203 fundamental shift in architecture towards *parallelization*, just as the shift to the *relational model* was a one-time
204 shift. As *relational databases* evolved to greater efficiencies over decades, so too will *big data* technologies
205 continue to evolve. Many of the conceptual underpinnings of *big data* have been around for years, but years
206 since the mid 2000's have seen an explosion in scaling technologies and their maturation and application to
207 scaled data systems.

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

214 The guidance for the choice of *big data* architectures is driven by four *big data* characteristics.

215 **4.1 Volume**

216 Volume is one of the characteristics of datasets that is most associated with *big data*. Volume represents the
217 extensive amount of data available for analysis to extract valuable information. The assumption that you can
218 extract the most value by analysing as much of the volume of data as possible was one of the primary drivers
219 for the creation of the new scaling technologies.

220 **4.2 Velocity**

221 Velocity is the rate of flow at which the data is created, stored, analysed or visualized. *Big data* velocity means
222 a large quantity of data needs to be processed in a short amount of time. Dealing with high velocity data is
223 commonly referred to as techniques for *streaming data*.

224 **4.3 Variety**

225 Variety represents the need to analyse data from a number of domains and a number of *data types*. The
226 variety of data was handled through transformations or pre-analytics to extract features that would allow
227 integration with other data. The wider range of data formats, logical models, timescales, and semantics, which
228 is desirable to use in analytics, complicates the integration of the variety of data. *Metadata* is increasingly used
229 to aid in the integration

230 **4.4 Variability**

231 Variability refers to changes in data rate, format/structure, semantics, and/or quality that impact the supported
232 application, analytic, or problem. Impacts can include the need to refactor architectures, interfaces,
233 processing/algorithms, integration/fusion, storage, applicability, or use of the data.

234

235
236
237
238

Annex A (informative)

Additional Related Concepts

239 The big data vocabulary consists of a number of additional terms and concepts that are in common usage.

240 **A.1 Volatility**

241 Volatility in Big Data usage refers to *data volatility*, or the rate of change of data over time

242 **A.2 Distributed File System**

243 Given the use of *horizontal scaling*, one new concept is a named set of records distributed across nodes of a
244 *cluster* treated as a unit.

245 **A.3 Non-relational Models**

246 Logical data models that do not follow *relational algebra* for the storage and manipulation of data

247 **A.4 Non-relational Database**

248 A non-relational database is a database that does not follow a *relational model*. NoSQL, which is typically
249 translated as “no-SQL” or “not only SQL”, is the term in common usage to refer to databases that do not
250 conform to a *relational model*.

251 **A.5 Scatter-gather**

252 Processing large data sets that are distributed across nodes of a *cluster* requires an algorithmic change to do
253 alter the algorithm to process the data on each node, and then return the results to process across the entire
254 dataset. For example, MapReduce is an implementation of a scatter-gather process for data processing

255 **A.6 Computational Portability**

256 Computational portability is the vocabulary for the movement of the computation to the location of large
257 volume datasets, rather than moving a large dataset to a different processing location.

258 Note: related to the prior concept of agents

259 **A.7 Data Science**

260 Data Science is the extraction of actionable knowledge from data through a process of discovery, or
261 hypothesis and hypothesis testing.

262 **A.8 Value**

263 Value represents the benefit to the organization of the actionable knowledge derived from an analytic system.
264 This term is not new to *big data*, but is often ascribed to *big data* due to the understanding that data has
265 potential value that was typically not considered previously. Value arises as an output of the implementation
266 of *big data* systems.

267 **A.9 Veracity**

268 Veracity refers to the completeness and accuracy of the data and relates to the data quality issues in
269 existence for a long time. If the analytics are causal, then the quality of every data element is extremely
270 important. If the analytics are correlations or trending over massive volume datasets, then individual bad
271 elements could be lost in the overall counts and the trend will still be accurate.

272 **A.10 Unstructured Data**

273 Unstructured data, such as text, image, video, and relationship data, have been increasing in both volume and
274 prominence. While modern *relational databases* tend to have support for these types of data elements, their
275 ability to directly analyse, index, and process them has tended to be both limited and accessed via non-
276 standard SQL extensions. The need to analyse unstructured or semi-structured data has been present for
277 many years. However, the *big data* paradigm shift has increased the emphasis on the value of unstructured or
278 relationship data, and also on different engineering methods that can handle data more efficiently.

Annex B (informative)

Cross-Cutting Concepts of Big Data

279
280
281
282

283 The development of Big Data systems has implications for a number of technological areas of discussion and
284 standardization. In this section we discuss the relationships to other areas of investigation.

285 **B.1 Metadata**

286 *Metadata* is data about data or data elements, including the description of the processing history of the data.
287 As Big Data systems are architected to perform *distributed data processing* including data that is external and
288 not under the control of the big data system, the use of *metadata* becomes an increasingly important concept.
289 As open data is reused for purposes far removed from its collection, it is important that *metadata* be
290 associated with any data that is made available to others.

291 **B.2 Analytics/Statistics**

292 The development of algorithms for the analysis of data previously did not consider the requirements of
293 *distributed data processing*, the data was typically held locally. For *big data*, algorithms must be adapted to
294 explicitly accommodate the particular distribution of data across nodes.

295 **B.3 Cluster Computing**

296 *Big data* systems can use the coordination of nodes in a cluster to achieve scalability in the *distributed data*
297 *processing*. The nodes can be a physical compute system leveraging software to cause the nodes to function
298 as a unit, or it can insert a virtualized interface on top of the physical systems to realize the advantages of
299 cloud computing.

300 **B.4 Cloud Computing**

301 Cloud computing is a paradigm for enabling network access to a scalable and elastic pool of shareable
302 physical or virtual resources with self-service provisioning and administration on demand. There are several
303 key characteristics often present in for cloud computing deployments including: broad network access,
304 measured service, multi-tenancy, on-demand self-service, rapid elasticity and scalability, and resource
305 pooling. Cloud computing is an infrastructure model for the development of a Big Data system. To achieve the
306 needed scalability, systems can leverage infrastructure as a service (IaaS), data platform software can to
307 provide a platform as a service (PaaS), or an application can be provided as software as a service (SaaS)

308 **B.5 Security**

309 *Big data* systems have additional security concerns due to the distributed nature of the *data processing*.
310 Additional vulnerabilities arise for example with the distributed ownership and control of the physical compute
311 and network infrastructure, as well as the control across each level of the software and storage frameworks.

312 **B. 6 Privacy**

313 Privacy concerns are heightened now that so much data is available about individuals from the web, social
314 media, sensors, and so forth. Privacy in a wide sense of the term is information control by an individual, which
315 is not only to prevent data usage disadvantageous to the individual data subjects but also to facilitate data
316 usage beneficial to them. A typical issue in the former respect is that integration of datasets may create
317 personally identifiable information (PII), even if none of these datasets contains PII independently. On the
318 other hand, the major problem in the latter respect is that personal data is mostly owned by organizations, and
319 typically not by the data subjects themselves. This makes it hard for individuals to use their data for their own

320 sake, and for organizations to use personal data owned by other organizations, so that noone can utilize big
321 and deep personal data.

322 **B.7 SQL**

323 SQL (Structured Query Language) is a standard interactive and programming language designed for
324 querying, updating, and managing data and data set in the database management system. SQL was first
325 published as ISO International Standard (ISO/IEC 9075) in 1987, and it has been revised to include a larger
326 set of features as the query language for Information technology. SQL is designed for manipulating structured
327 data, and it is also fast becoming the default language for big data analytics. SQL provides a mature and
328 comprehensive framework for data access supporting a broad range of advanced analytical features.
329 Analytics is a must-have functional component of data warehouse and big data. Modern SQL databases
330 support the discovery of columns across a wide range of data set: not only relational table/views, but also
331 XML, JSON, spatial objects, image-style objects (Binary Large Objects and Character Large Objects), and
332 semantic objects.

333 **B.8 Parallel Computing**

334 Big Data typically refers to distributed data-intensive processing across the nodes of a *cluster*. The simulation
335 community has been developing methods for compute-intensive processing across large clusters of nodes for
336 many years. Given that both approaches represent extrema cases for large scale computation and large scale
337 data analysis, techniques from both will be leveraged for spectrum of capabilities needing both compute-
338 intensive and data-intensive computation.

339 **B.9 Internet of Things**

340 More and more data are being created, along with computing systems capable of analysing data. Users want
341 to leverage the amount of data available from a variety of sensors and other data generators. This provides
342 efficient predictive analytics to manage and control networked solutions. Typical technological advances in
343 sensors, and the deployment of IPV6 to provide Internet connectivity to sensors creates the need for a *big*
344 *data* system that can handle high velocity *streaming data* from a number of sources. This is in contrast to high
345 volume big data systems that typically run batch jobs over a relatively small number of large datasets. This
346 difference in the characteristics of the datasets has direct implications on the architecture and methods used
347 for data analysis.

348 **B.10 Programming Languages**

349 An analysis of extended data by using statistical computing is the fundamental approach for *big data*.
350 Customers can develop big data analytics system by using general-purpose programming languages. On the
351 other hand, R is a programming language and software environment for statistical computing and graphics for
352 Statistical Computing. The R language is widely used among statisticians and data miners for developing
353 statistical software and data analysis. R is freely available under the GNU General Public License since 1995,
354 and is adopted onto various operating systems. R is widely used for developing statistical computing
355 applications that analyse big data.

356 The needs for *distributed data processing* have led to a number of new programming and query languages
357 suited to the development of *big data* systems, as well as new processes. An example of new programming
358 languages and frameworks is the Spark framework for analytics. New data processing languages such as Pig
359 have also been developed. New query languages address data in NoSQL systems. New processes include
360 scatter-gather for *distributed data processing*.

361

362

Bibliography

- 363 [1] ISO 2382-1:1993, *Information technology – Vocabulary — Part 1: Fundamental terms*
- 364 [2] ISO 2382-3:1987, *Information processing systems - Vocabulary — Part 3: Equipment technology*
- 365 [3] ISO 2382-17:1999, *Information technology - Vocabulary — Part 17: Databases*
- 366 [4] ISO 2382-18:1999, *Information technology - Vocabulary — Part 18: Distributed data processing*
- 367 [5] ISO 19784-4:2011, *Information technology – Biometric application programming interface — Part 4:*
368 *Biometric sensor function provider interface*