



1 DDI Working Paper Series -- Best Practices, No. 5

2 **Subject**

3 Controlled Vocabularies (2009-02-22)

4 **Document identifier:**

5 <http://dx.doi.org/10.3886/DDIBestPractices05>

6

7

8 **Authors:**

9 Peter Granda, Stefan Kramer, Jenny Linnerud, Hans Jørgen Marker, Ken Miller,
10 Mary Vardigan

11 **Editors:**

12 Ken Miller

13 **Target audience:**

14 Communities that want to select, develop, or extend controlled vocabularies for use
15 with DDI.

16 **Abstract:**

17 The benefits of using controlled vocabularies within DDI metadata creation are the
18 primary underlying principles behind this Best Practice document.

19 Use of a controlled vocabulary enhances consistency and efficiency in the
20 production of DDI metadata, gives precision in searching the DDI metadata, and
21 allows semantic and technical interoperability between organisations creating DDI
22 instances.

23 **Status:**

24 This document is updated periodically on no particular schedule. Send comments to
25 editor: ddi-bp-editors@icpsr.umich.edu



26

27 **Table of Contents**

28

29	INTRODUCTION	3
30	1.1 Problem statement	3
31	1.2 Terminology	3
32	2 BEST PRACTICE SOLUTION	3
33	2.1 Definitions	3
34	2.2 Best Practice behavior	4
35	2.3 Discussion	6
36	2.4 Example	7
37	3 REFERENCES	10
38	3.1 Normative	11
39	APPENDIX A. ACKNOWLEDGMENTS	12
40	APPENDIX B. REVISION HISTORY	14
41	APPENDIX C. LEGAL NOTICES	15
42		



43

44 1 Introduction

45 1.1 Problem statement

46 Dependence on natural language alone in indexing and searching brings with it a number of
47 problems that use of a controlled vocabulary (CV) can address. For a fuller explanation of
48 the problems a CV addresses, see the Discussion section.

49 This Best Practice does not address creating CVs themselves. Within the DDI community,
50 the DDI Alliance Controlled Vocabularies Working Group is focusing on this task.

51 This best practice addresses:

- 52 • **Technical aspects:** CV publication format; application requirements for handling
53 hierarchies, if present in a CV, for indexing and searching; separation of index terms
54 and structural description in a CV; and management of the CV's content.
- 55 • **Business practices:** Selection of a platform for publishing CVs; decision about term
56 structure in CVs; relationship identification in CVs; versioning and cross-referencing
57 of terms in CVs; policies governing CV maintenance; mapping between different
58 CVs; documentation of CVs; and CVs for system use.

59 1.2 Terminology

60 The key words *must*, *must not*, *required*, *shall*, *shall not*, *should*, *should not*, *recommended*,
61 *may*, and *optional* in this document are to be interpreted as described in [RFC2119].

62 Additional DDI standard terminology and definitions are found in

63 <http://www.ddialliance.org/bp/definitions>

64 2 Best Practice Solution

65 2.1 Definitions

66 Genericode: Genericode defines a standard format for defining code lists (also known as
67 enumerations or controlled vocabularies).

68 Genericode aims to provide the following:

- 69 • A standard model and XML representation for the contents of a code list
- 70 • A standard model and XML representation for data associated with items in a code
71 list



Data Documentation Initiative

- 72 • A standard model and XML representation for how new code lists are derived from
73 existing code lists

74 Community: In this document, the term *community* is used to identify any grouping of
75 personal or organizational entities, at different levels of formal organization, that are
76 considering or undertaking implementation of DDI. Examples: a national statistical service,
77 a data producer, an archive, a consortium of data archives.

78 Controlled vocabulary (CV): Broadly speaking, a CV can range from a short list of clearly
79 defined, mutually exclusive, and exhaustive terms, which are the only choices for usage in a
80 specific context (e.g., populating certain DDI elements or attributes) through a classification
81 to something as complex as a thesaurus with thousands of terms and term relationships. A
82 CV has also been described as “A set of subject terms, and rules for their use in assigning
83 terms to materials for indexing and retrieval.”

84 (<http://www.cs.cornell.edu/wya/diglib/MS1999/Glossary.html>)

85 In a CV, a term consists of one or more words used to represent a concept (example: “fear”;
86 “females”; “child care”). Terms are selected from natural language for inclusion in a
87 controlled vocabulary.

88

89 Pre-coordinated vs. Post-coordinated systems in CVs:

- 90 • In pre-coordinated CV systems, multiple concepts are brought together in one term.
91 An illustrative example is the Library of Congress Subject Headings (LCSH), which
92 yield entries such as: “Insurance, Unemployment --Switzerland --Statistics.” This
93 method allows for disambiguation of the relationship of the concepts in the term that
94 might not be possible in post-coordinated systems, such as whether a term is a
95 qualifier of another.

- 96 • In post-coordinated or faceted systems, concepts are kept broad and separate and
97 selected and joined in the process of searching with Boolean operators.
98 A representation of the above LCSH in this system could be “Insurance AND
99 Statistics AND Switzerland AND Unemployment” – note that entry order in the query
100 has no relevance here. An example of such a system is the American Psychological
101 Association’s *Thesaurus of Psychological Index Terms*.

102

103 **2.2 Best Practice behavior**

104 **Selection or Creation**

- 105 • Research whether a CV already exists that maps well with the knowledge domain to
106 be covered by your activity, which could mean taking a subset of that existing CV.



Data Documentation Initiative

- 107 • Research whether the desired CV has already been published in a format suitable to
108 local application, such as Genericcode. (Recommend to creators of CVs to publish
109 the machine-actionable part of the CV in Genericcode.)

- 110 • In selecting an existing CV for use, choose one that is as large and complex as
111 necessary, but as small and simple as possible, to help optimize both indexing and
112 searching.

- 113 • Multiple languages should be considered in the initial creation or selection of a CV.

- 114 • If selecting or creating a complex CV, decide upfront whether it is or will be pre- or
115 post-coordinated. Post-coordination will allow for more flexibility and lend itself to
116 Boolean searching by end users.

- 117 • If you are merging two or more complex CVs, consider the implications if some of
118 them are pre- and some are post-coordinated.

- 119 • In making a CV publicly available, consider releasing it with the least possible
120 restrictions (e.g., Creative Commons or Public Domain).

- 121 • Multiple communities may join to create a CV derived from existing ones.

- 122 • If a new CV is created, the reason for this activity should be explained as part of its
123 documentation (see, for instance, ANSI/NISO Z39.19-2005 section 9.5 [see
124 References section]). If it is made publicly available, its creation should be
125 announced to the DDI community and any other relevant user communities.

- 126 • If terms are likely to be misunderstood by the CV's users, more detailed scope notes
127 should be provided for them.

- 128 • Handling of phrases: multi-word terms (e.g., "child care" or "traffic regulation") need
129 to stay bound as a phrase for handling by applications.

- 130 • If your CV contains descriptions of structure, these should be clearly distinguished
131 from terms to insure accuracy in indexing and processing. Descriptions of
132 hierarchies (e.g., "Age Groups" or "Vehicles by size") and other text that describes
133 the arrangement of the CV, but does not constitute CV terms, need to be excluded
134 from being able to be assigned as CV terms.

- 135 • The storage container for the CV should allow for: persistent identifiers, version
136 control and dating. This facilitates machine-actionability and human understanding
137 of the CV.



Data Documentation Initiative

- 138 • Consider use of CVs for system-to-system coordination, such as consistent access
139 control to shared resources among different communities.

140 **Maintenance, Use of, and Additions to CVs**

- 141 • Define what requires additions or revisions to the CV, such as: new organizations
142 contributing to or being affected by the CV; new data acquisitions; new services or
143 user communities; changes in the reality the CV describes; changes in language
144 usage or jargon; etc.

- 145 • If the community of CV users broadens over time, it may become necessary to
146 translate the CV into other languages.

- 147 • Whether using one or more complex CVs, the relationships between terms (e.g.,
148 “see” or “related terms” or “narrower term”) have to be defined and represented
149 consistently as the CV evolves.

- 150 • Deprecated terms in more complex CVs need to be identified as such, and refer to
151 their “use instead” term.

- 152 • Assign only the most specific CV term if a hierarchy of terms is exposed, so broader
153 terms can automatically be included in a query.

- 154 • Multiple communities that have agreed to create a CV derived from existing CVs
155 should create and maintain a mapping between the existing and derived CVs.

156 **2.3 Discussion**

157 Use of a controlled vocabulary in indexing, documenting, or searching for resources by use
158 of their assigned metadata will enhance: precision; consistency; temporal, spatial and
159 topical comparability; semantic and technical interoperability; multilingual access; efficiency;
160 and harmonization. These benefits of a controlled vocabulary are the primary underlying
161 principles behind this Best Practice document.

162 Legal considerations including applicable laws regarding copyright for using large parts of,
163 or entire, existing controlled vocabularies, should be taken into account. Detailed
164 discussion of legal issues is outside the scope of this Best Practice.

165 The need for any crosswalks between CVs that cover different domains of knowledge
166 should be considered, but is outside of the scope of this Best Practice.



167 **2.4 Example**

168 There are CVs embedded in the DDI standard already. These are short lists of clearly
169 defined, mutually exclusive, and exhaustive terms, which have vocabularies which are
170 considered fixed. Examples of these are:

171 **ValueTypeCodeType:** Indicates value type.

- 172 • *Restricts:* xs:NMTOKEN
- 173 • (pattern is: code, colon, term, dash, definition)
 - 174 ○ Code: GreaterThan - Greater Than
 - 175 ○ Code: LessThan - Less Than
 - 176 ○ Code: Equal - Equal
 - 177 ○ Code: GreaterThanOrEqual - Greater Than or Equal
 - 178 ○ Code: LessThanOrEqual - Less Than or Equal
 - 179 ○ Code: NotEqual - Not Equal

180 **AggregationMethodCodeType:** A list for describing aggregation methods.

- 181 • *Restricts:* xs:NMTOKEN
- 182 • (pattern is: code, colon, term, dash, definition)
 - 183 ○ Code: Sum - Sum
 - 184 ○ Code: Average - Average
 - 185 ○ Code: Count - Count
 - 186 ○ Code: Mode - Mode
 - 187 ○ Code: Median - Median
 - 188 ○ Code: Maximum - Maximum
 - 189 ○ Code: Minimum - Minimum
 - 190 ○ Code: Percent - Percent (Percentages are used to express how large
 - 191 one quantity is relative to another quantity)
 - 192 ○ Code: CumulativePercent - Cumulative percent (percentage of items
 - 193 in its frequency distribution which are equal to or lower than the
 - 194 current item [maximum value is 100%])
 - 195 ○ Code: PercentileRank - Percentile rank (The percentile rank of a item
 - 196 is the percentage of items in its frequency distribution which are lower
 - 197 [cannot reach 100%])

198 The DDI Controlled Vocabulary Group is looking at other areas of the DDI in which a CV
199 would enhance computer actionability, but the list of terms is not exhaustive, and the
200 terminology might change over time. These would be published in Genericode.

201 Examples of this kind of CV are:



Data Documentation Initiative

202 Suggested CV for Collection Mode (Data Collection Module)

- 203 • Interview
 - 204 ○ Face-to-face
 - 205 ○ Telephone
 - 206 ○ E-mail
 - 207 ○ CATI
 - 208 ○ CAPI
 - 209 ○ Mixed-mode
- 210 • Self-completed questionnaire:
 - 211 ○ Paper/pencil
 - 212 ○ Web-based
 - 213 ○ CASI
 - 214 ○ ACASI
- 215 • Coding
- 216 • Transcription
- 217 • Compilation
- 218 • Synthesis
- 219 • Recording
- 220 • Simulation
- 221 • Observation
- 222 • Experiments
- 223 • Focus Group
- 224 • Other
- 225

226 Suggested CV for attribute "Role":

- 227 • Data Collector
- 228 • Data producer
- 229 • Depositor
- 230 • Distributor
- 231 • Copyright holder
- 232 • Research Investigator
- 233 • Other
- 234

235 Within the Council of European Social Science Data Archives (CESSDA) further CVs are
236 being considered to enhance searching and additional proposed harmonization and
237 comparability services that would be available via the CESSDA portal.

238 Examples here are:

239 (Pattern is: hierarchy level, dashes indicating level, term)

240 01 - - SAMPLING PROCEDURES



Data Documentation Initiative

- 241 02 - - - - COMPLETE COUNT
- 242 02 - - - - PROBABILITY SAMPLE
- 243 03 - - - - - SIMPLE RANDOM SAMPLE
- 244 03 - - - - - SYSTEMATIC SAMPLE
- 245 03 - - - - - STRATIFIED SAMPLE
- 246 03 - - - - - CLUSTER SAMPLE
- 247 03 - - - - - MULTIPHASE SAMPLE
- 248 03 - - - - - TIME SAMPLE
- 249 02 - - - - NONPROBABILITY SAMPLE
- 250 03 - - - - - PURPOSIVE SAMPLE
- 251 04 - - - - - - QUOTA SAMPLE
- 252 04 - - - - - - RANDOM WALK SAMPLE
- 253 03 - - - - - ACCIDENTAL SAMPLE
- 254 04 - - - - - - VOLUNTEER SAMPLE
- 255 04 - - - - - - CONVENIENCE SAMPLE
- 256 01 - - DATA FORMAT
- 257 02 - - - - TEXTUAL DATA
- 258 02 - - - - NUMERIC DATA
- 259 02 - - - - ALPHA/NUMERIC DATA
- 260 02 - - - - IMAGE DATA
- 261
- 262 01 - - DATA STRUCTURE
- 263 02 - - - - RECTANGULAR DATA
- 264 02 - - - - RELATIONAL DATA



Data Documentation Initiative

265 02 - - - HIERARCHICAL DATA

266

267 Example of the development of a community-based thesaurus:

268 CESSDA has already constructed a CV to act as a common topical classification scheme to
269 describe the overall coverage of the resources held at individual member archives.

270 The process involved consideration of existing published classification schemes, such as
271 the Library of Congress Subject Headings. Several CESSDA member archives had created
272 CVs themselves and through comparing these it was agreed that a new common CV would
273 not be that difficult to create.

274 After the common structure had been agreed each member archive had the choice to either
275 adopt the new CV or map their existing one to it. The terminology of about 140 terms was
276 then translated into the languages of the different member archives.

277 When it came to the adoption of a common thesaurus for allocation of keywords and
278 concepts CESSDA chose the existing UKDA HASSET thesaurus. This was reduced and
279 restructured slightly to remove UK specific concepts and align it more to a European
280 resource. The terms (over 4,000) have now been translated into 9 of the 21 languages of
281 the present CESSDA membership.

282 **3 References**

283 Genericcode <http://www.genericcode.org/>

284 DDI Technical Specification, April 2008, part I, section 4.3: *Controlled Vocabularies*.
285 Available via: http://sourceforge.net/projects/ddi-alliance/files/Data%20Documentation%20Initiative/DDI%203.0%20%282008-04-28%29/DDI_3_0_2008-04-28_Documentation_XMLSchema.zip/download

286 Best Practice on Versioning and Publication: <http://dx.doi.org/10.3886/DDIBestPractices08>
287

288 DDI Alliance Controlled Vocabularies Working Group: <http://www.ddialliance.org/alliance/working-groups>
289

290 Guidelines for the Construction, Format, and Management of Monolingual Controlled
291 Vocabularies: <http://www.niso.org/standards/resources/Z39-19-2005.pdf>

292 ANSI/NISO Z39.19-2005: Guidelines for the Construction, Format, and Management of
293 Monolingual Controlled Vocabularies:
294 http://www.techstreet.com/cgi-bin/detail?product_id=1262086



Data Documentation Initiative

295 ISO 5964: Guidelines for the establishment and development of multilingual thesauri:
296 http://www.iso.org/iso/iso_catalogue/catalogue_tc/catalogue_detail.htm?csnumber=12159

297 Creative Commons: <http://creativecommons.org/>

298 **3.1 Normative**

299

300 [RFC2119] S. Bradner, Key words for use in RFCs to Indicate Requirement Levels,
301 <http://www.ietf.org/rfc/rfc2119.txt>, IETF RFC 2119, March 1997.

302

303 OASIS, Best Practice, [http://www.oasis-open.org/committees/uddi-](http://www.oasis-open.org/committees/uddi-spec/doc/bp/uddi-spec-tc-bp-template.doc)
304 [spec/doc/bp/uddi-spec-tc-bp-template.doc](http://www.oasis-open.org/committees/uddi-spec/doc/bp/uddi-spec-tc-bp-template.doc), 2003



305

306 **Appendix A. Acknowledgments**

307 The following individuals were members of the DDI Expert Workshop held 10-14 November
308 2008 at Schloss Dagstuhl, Leibniz Center for Informatics, in Wadern, Germany.

309 Nikos Askitas, Institute for the Study of Labor (IZA)

310 Karl Dinkelmann, University of Michigan

311 Michelle Edwards, University of Guelph

312 Janet Eisenhauer, University of Wisconsin

313 Jane Fry, Carleton University

314 Peter Granda, Inter-university Consortium for Political and Social Research (ICPSR)

315 Arofan Gregory, Open Data Foundation

316 Rob Grim, Tilburg University

317 Pascal Heus, Open Data Foundation

318 Maarten Hoogerwerf, Data Archiving and Networked Services (DANS)

319 Chuck Humphrey, University of Alberta

320 Jeremy Iverson, Algenta Technology

321 Jannik Vestergaard Jensen, Danish Data Archive (DDA)

322 Kirstine Kolsrud, Norwegian Social Science Data Services (NSD)

323 Stefan Kramer, Yale University

324 Jenny Linnerud, Statistics Norway

325 Hans Jørgen Marker, Danish Data Archive (DDA)

326 Ken Miller, United Kingdom Data Archive (UKDA)

327 Meinhard Moschner, GESIS - Leibniz Institute for the Social Sciences

328 Ron Nakao, Stanford University



Data Documentation Initiative

- 329 Sigbjørn Revheim, Norwegian Social Science Data Services (NSD)
- 330 Wendy Thomas, University of Minnesota
- 331 Mary Vardigan, Inter-university Consortium for Political and Social Research (ICPSR)
- 332 Joachim Wackerow, GESIS - Leibniz Institute for the Social Sciences
- 333 Wolfgang Zenk-Möltgen, GESIS - Leibniz Institute for the Social Sciences



334

335 **Appendix B. Revision History**

336

Rev	Date	By Whom	What
0.1	14/11/2008	Stefan Kramer	Final version from Dagstuhl
0.2	30/11/2008	Ken Miller	Abstract added and minor text and format changes.
0.3	2009-02-15	Stefan Kramer	Formatting changes (standard across all BPs) and reference link additions.

337



338

339 **Appendix C. Legal Notices**

340 Copyright © DDI Alliance 2009, *All Rights Reserved*

341

342 <http://www.ddialliance.org/>

343

344 Content of this document is licensed under a Creative Commons License:

345 Attribution-Noncommercial-Share Alike 3.0 United States

346

347 This is a human-readable summary of the Legal Code (the full license).

348 <http://creativecommons.org/licenses/by-nc-sa/3.0/us/>

349

350 You are free:

351

- to Share - to copy, distribute, display, and perform the work

352

- to Remix - to make derivative works

353

354 Under the following conditions:

355

- Attribution. You must attribute the work in the manner specified by the author or licensor (but not in any way that suggests that they endorse you or your use of the work).

356

357

- Noncommercial. You may not use this work for commercial purposes.

358

359

- Share Alike. If you alter, transform, or build upon this work, you may distribute the resulting work only under the same or similar license to this one. For any reuse or distribution, you must make clear to others the license terms of this work. The best way to do this is with a link to this Web page.

360

361

362

- Any of the above conditions can be waived if you get permission from the copyright holder.

363

364

- Apart from the remix rights granted under this license, nothing in this license impairs or restricts the author's moral rights.

365

366

367

368 **Disclaimer**

369

370

The Commons Deed is not a license. It is simply a handy reference for understanding the Legal Code (the full license) — it is a human-readable expression of some of its key terms. Think of it as the user-friendly interface to the Legal Code beneath. This Deed itself has no legal value, and its contents do not appear in the actual license.

371

372

373

374

375

Creative Commons is not a law firm and does not provide legal services. Distributing of, displaying

376

of, or linking to this Commons Deed does not create an attorney-client relationship.

377

Your fair use and other rights are in no way affected by the above.

378

379

Legal Code:

380

<http://creativecommons.org/licenses/by-nc-sa/3.0/us/legalcode>