



## 1 DDI Working Paper Series -- Best Practices, No. 4

### 2 **Subject:**

3 Workflows for Metadata Creation Regarding Recoding, Aggregation and Other Data  
4 Processing Activities (2009-03-21)

### 5 **Document identifier:**

6 <http://dx.doi.org/10.3886/DDIBestPractices04>

### 10 **Target audience:**

11 Producers of DDI metadata; researchers in the sense of both data collectors and  
12 analysts; data processors whenever and wherever processing occurs (data  
13 collection, data cleaning, dataset creation, archiving process, or other)

### 14 **Authors:**

15 Hans Jørgen Marker, Wolfgang Zenk-Möltgen, Wendy Thomas, Achim Wackerow

### 16 **Editors:**

17 Hans Jørgen Marker, Wendy Thomas

### 18 **Abstract:**

19 This best practice discusses the capturing, in DDI metadata, of the processes of  
20 data aggregation, recoding and data processing.

### 21 **Status: Draft**

22 This document is updated periodically on no particular schedule. Send comments to  
23 [ddi-bp-editors@icpsr.umich.edu](mailto:ddi-bp-editors@icpsr.umich.edu)



24	<b>INTRODUCTION.....</b>	<b>3</b>
25	1.1 Problem statement .....	3
26	1.2 Terminology .....	3
27	<b>2 BEST PRACTICE SOLUTION .....</b>	<b>4</b>
28	2.1 Definitions .....	4
29	2.2 Best Practice behavior .....	4
30	2.3 Discussion .....	7
31	2.4 Example .....	8
32	<b>3 REFERENCES .....</b>	<b>10</b>
33	3.1 Normative .....	10
34	<b>APPENDIX A. ACKNOWLEDGMENTS .....</b>	<b>11</b>
35	<b>APPENDIX B. REVISION HISTORY .....</b>	<b>13</b>
36	<b>APPENDIX C. LEGAL NOTICES .....</b>	<b>14</b>



37

## 38 **Introduction**

39 DDI 3 facilitates the creation of metadata at a variety of starting points from the hypothesis  
40 for a study through the capturing of legacy metadata. How and where one starts capturing  
41 metadata depends upon the data being described, the application within which it is used,  
42 and the organizational needs of the creators. The best practices on workflow provide  
43 guidelines for setting up metadata creation processes within different environments,  
44 identifying organizational and application features that impact the process structure,  
45 addressing salient questions/issues in setting up the process, and determining the  
46 implications of various starting points and process orders:

- 47 1. Metadata Creation Regarding Recoding, Aggregation, and Other Data Processing  
48 Activities (this document)
- 49 2. Archival Ingest and Metadata Enhancement [see References section]
- 50 3. Dissemination and Discovery: User Perspective [see References section]

### 51 **1.1 Problem statement**

52 Data transformations related to aggregation, recoding, and data processing need to be  
53 captured in the metadata to support both data processing (machine-actionable) and user  
54 understanding of the resulting data content. This information must support both reuse of  
55 metadata and data but also the review and evaluation of resulting research. Stakeholders  
56 lack information if data transformations are not documented. This kind of documentation is  
57 essential to tracking the evolution of a dataset over time, as indicated in the OAIS reference  
58 model that defines requirements for digital preservation. Stakeholders include: researchers,  
59 research councils/funding agencies, data producers, archivists, librarians, data and  
60 metadata users, registry managers, and research analysts.

### 61 **1.2 Terminology**

62 The key words *must*, *must not*, *required*, *shall*, *shall not*, *should*, *should not*, *recommended*,  
63 *may*, and *optional* in this document are to be interpreted as described in [RFC2119].  
64 Additional DDI standard terminology and definitions are found in  
65 <http://www.ddialliance.org/bp/definitions>



66

## 67 **2 Best Practice Solution**

### 68 **2.1 Definitions**

69 Open Archival Information System (OAIS): A reference model of the space community that  
70 governs general archival activities and policies. Includes:

71 SIP: Submission Information Package

72 AIP: Archival Information Package

73 DIP: Dissemination Information Package

74 Control operations: Methods to facilitate data control performed by the primary investigator  
75 or by the data archive. Specify any special programs used for such operations. The  
76 "agency" attribute maybe used to refer to the agency that performed the control operation.

77 Cleaning operations: Methods used to "clean" the data collection, e.g., consistency  
78 checking, wildcode checking, etc. The "agency" attribute permits specification of the agency  
79 doing the data cleaning.

80 Weighting: The use of sampling procedures may make it necessary to apply weights to  
81 produce accurate statistical results. Describe here the criteria for using weights in analysis  
82 of a collection. If a weighting formula or coefficient was developed, provide this formula,  
83 define its elements, and indicate how the formula is applied to data.

84 Logical record: A reference to a data record that is independent of its physical location. It  
85 may be physically stored in two or more locations.

86 NCubes: Describe the logical structure of an n-dimensional array, in which each coordinate  
87 intersects with every other dimension at a single point. The NCube has been designed for  
88 use in the markup of aggregate data.

89 Resource package: A resource package is a means of packaging any maintainable set of  
90 DDI metadata for referencing as part of a study unit or group. A resource package  
91 structures materials for publication that are intended to be reused by multiple studies,  
92 projects, or communities of users. A resource package uses the group module with an  
93 alternative top-level element called Resource Package that is used to describe maintainable  
94 modules or schemes that may be used by multiple study units outside of a group structure.

### 95 **2.2 Best Practice behavior**

96 Data processing is captured by a set of elements that allow both descriptive and machine-  
97 actionable content regarding control operations, cleaning operations, weighting, data



## Data Documentation Initiative

98 appraisal, and variable generation. While commonly occurring between data collection and  
99 production of the dataset, in fact, data processing can and does take place at many points  
100 along the data life cycle. Activities or events that can trigger the capture of metadata related  
101 to data processing may include the following (not an exhaustive list):

- 102 • Instrument development
- 103 • Confidentiality checks
- 104 • Transfer from data collection to data file
- 105 • Error correction
- 106 • Incorporating externally collected data items
- 107 • Creating indicator variables
- 108 • Harmonization
- 109 • Building classifications
- 110 • Comparison
- 111 • Creating aggregated statistics
- 112 • Preparation for analysis
- 113 • Creating new logical records
- 114 • Distribution of data
- 115 • Archival ingest and data management
- 116 • Creating subsets
- 117 • Linking datasets and/or data integration
- 118 • Complex file preparation
- 119 • Analysis

120 In effect, data processing is not a step in an overall process but occurs at many stages and  
121 is part of many other workflows throughout the overall process. Because of the ability of this  
122 metadata to capture specific machine-actionable commands, it can be used to process not  
123 only the data, but also automate the creation of metadata later in the process.



## Data Documentation Initiative

124 This best practice addresses the general process of data transformation and then identifies  
125 specific cases that may require special attention.

### 126 **Processing Event**

127 When a data processing event has occurred, new metadata should be created to reflect that  
128 event. This includes researchers and analysts who recode or restructure the content of the  
129 data in the process of analysis. The specific recoding or processing action should be  
130 captured from the statistical package or other processing tool and entered into the DDI  
131 format so that their results can be understood and evaluated by peers. For analysis this  
132 includes preparation of the data, the statistical methods applied, and the results. Processing  
133 metadata should be captured in a machine-actionable way (e.g., SAS or SPSS code,  
134 mathematical formula, or system-independent formal language) when possible, but should  
135 always include a human-readable explanation of the process. In addition to the specific  
136 processing commands, metadata should include the *agent* who executed the process,  
137 *when* it was done (specifically and in relation to the overall life cycle), the *purpose* of the  
138 process, and the *rationale* for use of the specific method. In addition to entering this  
139 information in the data processing elements, major events should also be noted in the Life  
140 cycle list in the Archive module. This makes it easier for both those managing the  
141 development of data and users to understand where and when data processing occurred  
142 throughout the workflow.

143 The DDI Processing Event element is a packaging mechanism that contains discrete  
144 events. Each event should contain the appropriate descriptive type, e.g., Control Operation  
145 and its associated Coding.

146 Control Operation and Cleaning Operation should be packaged with any related codes or  
147 data appraisal events that occurred as a single integrated process.

148 Weighting describes the process of determining overall or specific weights. When the  
149 weighting process results in one or more weight variables, the processing event should  
150 include both the weighting description and the generation code for each weight variable. If a  
151 study contains a standard weight, this should be entered as a numeric value in weighting. It  
152 is to be expected that future versions of the standard will have a machine-actionable home  
153 for the standard weight.

154 Data Appraisal Information contains separate elements for response rate, sampling error  
155 description, and other appraisal information. This should be coupled with any specific  
156 coding information related to the process used for determining sampling error or other  
157 appraisal. For surveys it is considered a best practice to include both the response rate and  
158 sampling error in all documentation.



## Data Documentation Initiative

159 Coding contains two structural types, General Instruction and Generation Instruction.  
160 General instruction captures processes that were used on large sections of data such as  
161 imputation processes, suppression rules, and handling of non-response to questions.  
162 General instructions express overall processing of the resulting data file. A general  
163 instruction for a subset of the data may override another general instruction for the whole  
164 file. Generation instructions provide processing information for specific variables. In addition  
165 to providing specific information to the user regarding the generation of the final dataset,  
166 these fields in conjunction with information on concepts and questions can be used to  
167 generate metadata in the logical product, physical data structure, and physical instance.  
168 When creating aggregate data expressed as NCubes, the aggregation process for creating  
169 specific cell contents is captured in aggregation within generation instruction.

170 Statistical summary data housed in physical instance does not require processing codes as  
171 these are standard weighted and unweighted frequencies with or without filters. Statistics  
172 created through the NCube creation process should be recorded with processing events.  
173 This would include the recoding of variables to created code representations such as age  
174 cohorts. In addition, NCubes require the creation of variables to represent specific  
175 measures such as counts and percents. These are captured in generation instruction to  
176 clarify how a count was defined or in the case of a percent exactly what is being used as the  
177 numerator and denominator in each table.

### 178 **2.3 Discussion**

179 The advantage to documenting data transformation is that the structure provides a  
180 systematic and coherent means of capturing a variety of processing events that can take  
181 place at many points along the life cycle. It facilitates access to the processing  
182 documentation and makes it available for reuse within and between studies.

183 The use of this best practice supports the use of data processing content to drive the  
184 production of data and metadata.

185 It packages the description of a process with the coding so that it can be referenced as a  
186 specific step from many points in the life cycle. It provides transparency for each step within  
187 a workflow, thereby aiding both understanding of how the data arrived at its current state at  
188 any point in the process workflow, and management of the process itself.

189 Capturing this type of metadata requires a change in work processes to capture the specific  
190 processing done within statistical software or other applications. Currently, much of this  
191 information is lost as it is not routinely maintained by the application along with the resulting  
192 data.

193 In addition, this type of workflow change provides the opportunity to capture the reasoning  
194 behind specific data processing activities.



## Data Documentation Initiative

195 Capturing this information can easily be perceived as an additional burden especially during  
196 the analysis process, which may have short term goals. The payoff for the researcher or  
197 group of researchers is improved quality of the result. It supports the ability to replicate the  
198 analysis for evaluation and validation. This approach supports best practices as they have  
199 been defined for scientific activities.

200 However, there are more easily identifiable payoffs for the data producer and data manager.  
201 Capturing the process commands in a single location allows them to be reused by  
202 reference. Processing event metadata can be used to drive applications for data and  
203 metadata creation. Aggregation processes can be captured and used to produce  
204 aggregated data from microdata on demand, as opposed to creating and storing it as a  
205 dataset. Put another way, the metadata can drive rule-based aggregate data presentation.

206 For production and management organizations that want to ensure consistent processing  
207 such as cleaning practices, confidentiality processes, etc., the processing event information  
208 and coding structures can be published and maintained as a resource package for reuse  
209 within the organization. See the Schemes Best Practice [see References section] for more  
210 information on preparing resource packages for reuse.

211 The definition of what constitutes an individual processing event that results in the  
212 transformation of data from one stage to another cannot be clearly defined. A process may  
213 result in a number of transient versions of the data between declared stages. The contents  
214 of a single Processing Event should be coherent (e.g., contain a single type of processing  
215 event and its related coding) and sufficiently discrete to support clear referencing and  
216 execution of the command codes. When does a process have a stage one and a stage two  
217 and what constitutes the processing event? Is it each individual step (e.g., code) or a  
218 combination of a number of steps? If the steps within a process must occur in a specified  
219 order, then they should be listed as separate process events.

220 Currently there is no machine-actionable means of providing sequential information for the  
221 processing of coding; however, the order in which steps were taken in an actual chain of  
222 events can be captured in the Life cycle list.

### 223 **2.4 Example**

224 This example shows the use of Processing Event as packaging for a single step (e.g.,  
225 cleaning instructions and coding).

```
226 <d:ProcessingEvent isIdentifiable="true" id="PE_1" isDerived="false">  
227 <d:CleaningOperation>  
228 <r:Description>For all NCubes without suppression and whose contents are  
229 additive. The contents of the discrete cells in the NCube are aggregated  
230 and compared to the universe value. All errors have been corrected or  
231 noted.</r:Description>  
232 </d:CleaningOperation>
```





## Data Documentation Initiative

```
233 <d:CleaningOperation>
234 <r:Description>For all geographic levels that roll up into a parent level
235 (e.g., U.S. Counties within U.S. States). All cells where suppression is
236 not used and whose contents are additive, have been aggregated and
237 compared to the value in the parent geography. All errors have been
238 corrected or noted.</r:Description>
239 </d:CleaningOperation>
240 <d:Coding isIdentifiable="true" id="Coding_1">
241 <d:GenerationInstruction>
242 <d:SourceVariable
243 isReference="true"><r:ID>V1</r:ID><d:Mnemonic>BaseAge</d:Mnemonic>
244 <r:Description>Recodes discrete age into 3 age cohorts; under 18, 18 to
245 64, and 65 and over.</r:Description>
246 <r:Command>
247 <r:CommandText formalLanguage="SPSS">RECODE BaseAge (Lowest thru 17=1) (18
248 thru 64=2) (65 thru Highest=3) INTO AGE.EXECUTE.</r:CommandText>
249 </r:Command>
250 </d:Coding>
251 </d:ProcessingEvent>
252
253
254 Representation/CodingInstructionReference in Variable "AGE" will reference
255 "Coding_1"
256
257 Parallel of Processing Event content and listing in life cycle

258 <r:LifecycleInformation>
259 <r:LifecycleEvent isIdentifiable="true" id="le_1">
260 <r:EventType>Cleaning Operation</r:EventType>
261 <r>Date><r:SimpleDate>2008-04-01</r:SimpleDate></r>Date>
262 <r:AgencyOrganizationReference isReference="true">
263 <r:ID>mpc</r:ID>
264 </r:AgencyOrganizationReference>
265 <r:Description>Completed cleaning operations verifying cell counts within
266 NCubes and cell counts within geographies which roll completely into a
267 higher level for which they provide comprehensive and complete
268 coverage..</r:Description>
269
```



270

271 **3 References**

272 DDI Best Practice: Workflows - Archival Ingest and Metadata Enhancement:

273 <http://dx.doi.org/10.3886/DDIBestPractices03>

274

275 DDI Best Practice: Workflows - Data Discovery and Dissemination: User Perspective:

276 <http://dx.doi.org/10.3886/DDIBestPractices02>

277

278 DDI Best Practice: DDI 3.0 Schemes:

279 <http://dx.doi.org/10.3886/DDIBestPractices07>

280 **3.1 Normative**

281

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

284

285 OASIS, Best Practice, [http://www.oasis-open.org/committees/uddi-spec/doc/bp/uddi-  
spec-tc-bp-template.doc](http://www.oasis-open.org/committees/uddi-spec/doc/bp/uddi-<br/>286 spec-tc-bp-template.doc), 2003

287

288 **Appendix A. Acknowledgments**

289 The following individuals were members of the DDI Expert Workshop held 10-14 November  
290 2008 at Schloss Dagstuhl, Leibniz Center for Informatics, in Wadern, Germany.

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

292 Karl Dinkelmann, University of Michigan

293 Michelle Edwards, University of Guelph

294 Janet Eisenhauer, University of Wisconsin

295 Jane Fry, Carleton University

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

297 Arofan Gregory, Open Data Foundation

298 Rob Grim, Tilburg University

299 Pascal Heus, Open Data Foundation

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

301 Chuck Humphrey, University of Alberta

302 Jeremy Iverson, Algenta Technology

303 Jannik Vestergaard Jensen, Danish Data Archive (DDA)

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

305 Stefan Kramer, Yale University

306 Jenny Linnerud, Statistics Norway

307 Hans Jørgen Marker, Danish Data Archive (DDA)

308 Ken Miller, United Kingdom Data Archive (UKDA)

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

310 Ron Nakao, Stanford University



Data Documentation Initiative

- 311 Sigbjørn Revheim, Norwegian Social Science Data Services (NSD)
- 312 Wendy Thomas, University of Minnesota
- 313 Mary Vardigan, Inter-university Consortium for Political and Social Research (ICPSR)
- 314 Joachim Wackerow, GESIS - Leibniz Institute for the Social Sciences
- 315 Wolfgang Zenk-Möltgen, GESIS - Leibniz Institute for the Social Sciences



316

317 **Appendix B. Revision History**

318

Rev	Date	By Whom	What
0.9	2008-02-08	Stefan Kramer	Removed date from filename to accommodate linking. Began revision history tracking.

319



320

## 321 **Appendix C. Legal Notices**

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

323

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

325

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

327 Attribution-Noncommercial-Share Alike 3.0 United States

328

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

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

331

332 You are free:

- 333 • to Share - to copy, distribute, display, and perform the work
- 334 • to Remix - to make derivative works

335

336 Under the following conditions:

- 337 • Attribution. You must attribute the work in the manner specified by the author or  
338 licensor (but not in any way that suggests that they endorse you or your use of  
339 the work).
- 340 • Noncommercial. You may not use this work for commercial purposes.
- 341 • Share Alike. If you alter, transform, or build upon this work, you may distribute  
342 the resulting work only under the same or similar license to this one. For any  
343 reuse or distribution, you must make clear to others the license terms of this  
344 work. The best way to do this is with a link to this Web page.
- 345 • Any of the above conditions can be waived if you get permission from the  
346 copyright holder.
- 347 • Apart from the remix rights granted under this license, nothing in this license  
348 impairs or restricts the author's moral rights.

349

350 **Disclaimer**

351

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

356

357 Creative Commons is not a law firm and does not provide legal services. Distributing of, displaying  
358 of, or linking to this Commons Deed does not create an attorney-client relationship.

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

360

361 **Legal Code:**

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