# Data Documentation Initiative (DDI) Technical Specification

## Part II:

## User Guide

**Version 3.0**

April 2008

*<ddi>* *Data Documentation Initiative*

*<ddi>* *Data Documentation Initiative*

48 # User Guide for DDI Version 3.0
49
50 Version 2:
51 Date: April 28, 2008
52 Wendy Thomas, Arofan Gregory, J Gager
53

54 # Table of Contents

# 118  1.0  Overall Structure

## 119  *1.1  Inline Content vs. Referenced Content*

120  The flexibility of a set of schemas means that you have a number of choices in
121  how to organize your DDI XML instance. Schema contents may be held
122  separately as maintained objects and then included by reference in a parent
123  document or included inline as a single document. Inclusion inline means that a
124  section is incorporated into the parent object in a hierarchical manner, with the
125  elements of each part listed in the main document. Inclusion by reference allows
126  you to create, say, a DDIInstance that references other DDIInstances which
127  contain its studyunit module(s), datacollection module(s), logicalproduct
128  module(s) etc. It is, in essence, a short file of external references, that when
129  resolved, result in the full document. Keep in mind that any decision you make
130  now does not preclude later disassembly of an inline document or assembly of  a
131  document containing external references into a single inline document.
132
133  Either approach or even combinations of the two are all permissible. The
134  approach taken depends on the anticipated use of the document and the needs
135  of the organization creating it. Issues that one should consider include:
136
137  Will any of the parts be used in multiple DDI XML instances?
138
139      For example, the data collection module of a census might be used in over
140      50 different logical products from public use micro-samples to aggregate
141      data tables. You may or may not want all of these logical products
142      expressed as a single document. If not, having the common parts
143      maintained as discrete objects allows you to use them in multiple
144      documents and still retain the relationships among the documents by their
145      mutual use of a common object.
146
147  Is it a simple study that is not anticipated to experience major changes, revisions,
148  or extensions?
149
150      Many of the studies found in archives are simple studies that are the result
151      of academic research. Such studies do not undergo extensive changes as
152      they are created, deposited, distributed, and analyzed.  In such cases, a
153      single inline document may be the easiest for the researcher to create.
154
155  In the final analysis, it is what is most expedient for the creator, publisher, or
156  archive.

## 157  *1.2  Top-Level Declarations*

158  All XML documents, whether using DTDs or Schemas, must declare their
159  structures at the head of the file. Because schema structures tend to use multiple

160  sources (multiple schemas), their declaration is both more complex and provides
161  for more options. In essence, the header provides the following information:
162
163  Declaration that this is an xml file:
164          <?xml version="1.0"?>
165
166  The primary schema's primary element tag:
167          Left angle bracket
168          abbreviation for the schema if used
169          colon
170          primary element tag
171          EXAMPLE:
172            <s:StudyUnit
173
174  The location of the primary schema including its URN filename and path (note
175  that the path can be to an internal copy of the schemas or an http path to a
176  remote copy):
177          xsi:schemaLocation
178          equal sign
179          open quote
180          schema URN
181          space
182          name of schema file including internal path or http path
183          close quote
184          EXAMPLE:
185            xsi:schemaLocation="ddi:studyunit:3_0 C:\\ddi\schemas\studyunit.xsd"
186
187  A namespace is the full URI of a schema or element. In the declaration, an
188  abbreviation or prefix is assigned to the XML schema namespace so that
189  elements from that schema can be uniquely identified with a [prefix]:[element
190  name] (see http://en.wikipedia.org/wiki/XML_Namespaces for additional
191  information and tutorial). All of the schemas required by the document must be
192  identified by namespace and assigned an abbreviation to be used by any
193  element found outside of the parent schema. Abbreviations are assigned to a
194  namespace by the "xmlns:" statement, which declares first the abbreviation and
195  then the namespace of the schema as ddi:<schema name>:version number. For
196  example, xmlns:r="ddi:reusable:3_0"
197          xmlns
198          colon
199          schema abbreviation
200          equal sign
201          open quote
202          schema URN or remote site
203          close quote
204          EXAMPLES:
205            xmlns:r="ddi:reusable:3_0"

206       xmlns:xhtml="http://www.w3.org/1999/xhtml"
207
208   All schemas except the top level instance.xsd have been assigned abbreviations
209   as a means of referencing them from within the schema definitions. Using these
210   abbreviations is a convenient convention if you want your documents to be more
211   easily recognizable by other human readers. You may alternately assign your
212   own and some XML editors assign default abbreviations. Systems will follow the
213   abbreviation pattern identified in the header. The following table provides the
214   abbreviation and URN used in the xmlns declarations, and the schema file name
215   used in the declaration of the primary schema.
216

| Abbr | URN ddi:schema:version | Schema name (ddi schemas) |
|---|---|---|
|  | ddi:instance:3_0 | instance.xsd |
| s | ddi:studyunit:3_0 | studyunit.xsd |
| d | ddi:datacollection:3_0 | datacollection.xsd |
| l | ddi:logicalproduct:3_0 | logicalproduct.xsd |
| c | ddi:conceptualcomponent:3_0 | conceptualcomponent.xsd |
| cm | ddi:comparative:3_0 | comparative.xsd |
| g | ddi:group:3_0 | group.xsd |
| pr | ddi:ddiprofile:3_0 | ddiprofile.xsd |
| a | ddi:archive:3_0 | archive.xsd |
| o | ddi:organizations:3_0 | organizations.xsd |
| p | ddi:physicaldataproduct:3_0 | physicaldataproduct.xsd |
| pi | ddi:physicalinstance:3_0 | physicalinstance.xsd |
| ds | ddi:dataset:3_0 | dataset.xsd |
| r | ddi:reusable:3_0 | reusable.xsd |
| dc | ddi:dcelements:3_0 | dcelements.xsd |
| xhtml | http://www.w3.org/1999/xhtml | [reference to standard] |
| xs | http://www.w3.org/XML/1998/namespace | xml.xsd |
| m1 | ddi:physicaldataproduct_ncube_normal:3_0 | physicaldataproduct_ncube_normal.xsd |
| m2 | ddi:physicaldataproduct_ncube_tabular:3_0 | physicaldataproduct_ncube_tabular.xsd |
| m3 | ddi:physicaldataproduct_ncube_inline:3_0 | physicaldataproduct_ncube_inline.xsd |
| m4 | ddi:physicaldataproduct_proprietary:3_0_Beta | physicaldataproduct_proprietary.xsd |

217
218   The first example below shows a document using the instance.xsd as the primary
219   schema. The primary element of instance.xsd is DDIInstance (parent of all other
220   elements in instance.xsd). In this example, the user has declared all the schemas
221   when in fact many are never used in the document. Note that the user has
222   provided the abbreviation "ns1" to the instance schema. This is common when
223   using XML editing software as they tend to assign prefixes for everything. With
224   the exception of this abbreviation, others were taken from the schema
225   declarations themselves. Both examples lack path information for the primary
226   schema, and so the xml document needs to reside in the same directory as the
227   schema. Once DDI 3.0 is officially published, we would normally use the official
228   maintained version of the schema or an agreed upon locally held copy,
229   accessible to the system processing the xml document.
230
231   <?xml version="1.0"?>
232   <ns1:DDIInstance

233    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
234    xsi:schemaLocation="ddi:instance:3_0 instance.xsd"
235    xmlns:ns1="ddi:instance:3_0"
236    xmlns:r="ddi:reusable:3_0"
237    xmlns:xhtml="http://www.w3.org/1999/xhtml"
238    xmlns:dc="ddi:dcelements:3_0"
239    xmlns:a="ddi:archive:3_0"
240    xmlns:g="ddi:group:3_0"
241    xmlns:cm="ddi:comparative:3_0"
242    xmlns:c="ddi:conceptualcomponent:3_0"
243    xmlns:d="ddi:datacollection:3_0"
244    xmlns:l="ddi:logicalproduct:3_0"
245    xmlns:p="ddi:physicaldataproduct:3_0"
246    xmlns:ds="ddi:dataset:3_0"
247    xmlns:pi="ddi:physicalinstance:3_0"
248    xmlns:m1="ddi:physicaldataproduct_ncube_normal:3_0"
249    xmlns:m2="ddi:physicaldataproduct_ncube_tabular:3_0"
250    xmlns:m3="ddi:physicaldataproduct_ncube_inline:3_0"
251    xmlns:m4="ddi:physicaldataproduct_proprietary:3_0_Beta"
252    xmlns:s="ddi:studyunit:3_0"
253    xmlns:pr="ddi:ddiprofile:3_0">
254
255         In the second example the user has declared the studyunit.xsd to be the
256    primary schema and used its abbreviation "s". The primary element in
257    studyunit.xsd is StudyUnit.
258
259    <?xml version="1.0"?>
260    <s:StudyUnit xmlns:s="ddi:studyunit:3_0"
261    xmlns:c="ddi:concept:3_0"
262    xmlns:l="ddi:logicalproduct:3_0"
263    xmlns:r="ddi:reusable:3_0"
264    xmlns:d="ddi:datacollection:3_0"
265    xmlns:xhtml="http://www.w3.org/1999/xhtml"
266    xmlns:xs="http://www.w3.org/2001/XMLSchema"
267    xsi:schemaLocation="ddi:studyunit:3_0 studyunit.xsd">
268
269    Note that the use of anything but the DDIInstance module (instance.xsd) as a
270    top-level element is reserved for non-public use of the DDI standard (e.g., within
271    a single organization during processing). Any publication of DDI for general use
272    should always use DDIInstance as the top-level element. This convention insures
273    that all external systems of a single standardized top level structure for all DDI
274    documents.
275
276

277 *1.3    Standard Model Contents*

| Complex Reusable Element | Description of Purpose and Usage |
|---|---|
| Identification is provided through use of MaintainableType as an extension base | This is the identification of the maintainable module. It must include an attribute id, as well as agency, version and versionDate. All elements within the module can inherit the Agency and Version from this complex element so it must be declared at the highest level. |
| Coverage | Coverage provides detail for Topical, Temporal, and Spatial Coverage. This should always be stated in the studyunit, and all sub-modules inherit this coverage unless they use this same element to constrain one or more of the coverage areas. A group should reflect the coverage of all studyunits included in the group. |
| OtherMaterial | OtherMaterial allows you to provide a citation, identifier, and type for any form of material (digital or otherwise) related to the study. This can be linked to any identifiable element in the document. This is available in all major modules so that OtherMaterial specific to a module can be placed at the most appropriate level. All OtherMaterial content is placed together near the top of a module. Storing OtherMaterial in the most relevant module facilitates keeping this information with the related identified objects if the instance is later broken down into its component parts. |
| Note | Note is a repeatable complex element located near the top of each module. Think of it in the same way you would a set of footnotes at the end of an article. The major features of Note are the ability to type the note and to attach it to any identifiable element in the module. A note can be created once and attached to as many elements as needed. Note was designed to be used during the production process to note questions, cleanups and verifications. Its single location and multiple linkages make it easy to edit or remove. |

278

## 1.3.1  Coverage

280  *1.3.1.1        Topical Coverage*
281  Topical coverage provides for both subject and keywords. It is an identifiable
282  object. Subject is intended to be a structured subject from a controlled
283  vocabulary such as MESH, LC, etc. Keywords are simple terms, but DDI allows

284  for the creation of a controlled vocabulary to further identify and relate keywords
285  between instances or with large collections.
286
287  *1.3.1.2 Temporal Coverage*
288  Temporal coverage provides information on the dates covered by the data within
289  the study unit or group. Each date is of the date structure described in section
290  2.7. At the top level (Study Unit or Group) Temporal coverage should provide
291  inclusive coverage. When used at a lower level, for example a single data file,
292  Temporal coverage can be used to constrain the coverage to a specific time or
293  time period within the defined top level coverage.
294
295
296  *1.3.1.3        Spatial (Geographic) Coverage*
297  Geographic coverage provides the following sections:
298    • A human-readable description of the geographic coverage which maps to
299       the Dublin Core geographic coverage element
300    • A bounding box (north, south latitudes and east, west longitudes) to
301       support geographic coordinate point search systems
302    • The lowest form of geographic description (point, line, polygon, or linear
303       ring)
304    • Spatial object
305    • Top Level and Lowest Level
306    • Summary Data Reference
307    • The relational structure of geographies
308    • Specific geographic locations
309
310  The human readable description should provide as much information on
311  geographic coverage as you would wish to see in a standard bibliographic
312  record. A bounding box consists of a North Latitude, South Latitude, West
313  Longitude, and East Longitude.  Longitude has a value in the range of -180
314  degrees to 180 degrees expressed as a decimal. Latitude has a value of -90
315  degrees to 90 degrees  expressed as a decimal. Bounding boxes provide a first
316  pass identification for coordinate value searches. While the area may seem quite
317  larger than the geography being covered (for example a bounding box of the
318  current United States includes most of Canada), in terms of the entire surface of
319  the earth, it is relatively specific. Internal bounding polygons can provide further
320  coordinate point detail. Geographic search system will also be able to make use
321  of Geographic Structure and Geographic Location details if they make use of
322  known gazetteer names and structure definitions.
323
324  The spatial object indicates whether the data is reported for points (a specific
325  address), a line (street, boarder), ring (area within X distance of a point in all
326  directions), or polygon (State, tract, place, block). This is not the level that the
327  data is collected at, but the level it is available in the dataset. This information
328  tells the geographer what he is able to do with the data in terms of aggregation.
329  Points can be aggregated to lines and polygons, but polygon level data cannot

330    be disaggregated to point data. Top level and Lowest level are the equivalents of
331    extent of coverage (for example Europe) and the smallest identifiable area (for
332    example Country). Both may be expressed as Names and as references to
333    specific levels in a GeographicStructure definition.
334
335    The last two are new features introduced in DDI 3.0 that allow for specific
336    identification of geographic structures (types of geographic areas like Countries,
337    States, Cities, etc.) and specific geographic locations (for example Germany,
338    France, Ghana, Japan, Canada, etc.). In addition a geographic polygon or
339    reference to shape files can be made at this level. These schemes can be large
340    and complex and the ability to create these once and reuse them via reference
341    was intended to both increase access to this level of geographic information, and
342    to provide a level of comparability by having multiple studies pointing to the same
343    geographic schemes. If this information is provided inline within the DDI Instance,
344    the information is held with GeographicStructureScheme and
345    GeographicLocationScheme found in conceptualcomponent.xsd. The appropriate
346    contents are then referenced by SpatialCoverage.
347
348    The purpose of the various geographic coverage descriptions is two-fold. First, it
349    provides information in a format that is more readily understood by geographers
350    and GIS systems. Second it provides access to geographic structure information
351    that was frequently only available by querying the data or by referring to outside
352    documents. Charts and tables providing coding information and relationships
353    generally required human interpretation in order to use them for accessing data.
354    DDI 3.0 attempts to provide this information in a form that can be processed for
355    both data discovery and data exploration. The major components and their
356    functions are listed below. It is important to have a sense of how this information
357    is used by others in order to understand the implications of excluding it from your
358    document.
359
360

| Element | Usage |
|---------|-------|
| Bounding Box | The bounding box is the north/south latitudes and east/west longitudes that bound the fullest extent of the geography being covered by the study. While this frequently encompasses more than the intended area (for instance the bounding box for the United States ends up including most of Canada), it is used as a first pass locator, by systems that search the world in terms of a coordinate point. Even a large, overextended box is only a small percentage of the earth's surface and therefore effectively limits the search area for this type of system. |

| Description | This is the human-readable description of the geographic coverage that maps to the Dublin Core "spatial coverage." It should be as informative as you wish, noting at minimum the extent of the coverage (top level or what the bounding box is bounding) and the smallest level of identifiable geography available (lowest level). Note that while DDI allows structure in this field through the use of xhtml tags, these tags will be lost if the information is transferred to Dublin Core. Make sure content is understandable without the structure elements if you anticipate using this field to populate a Dublin Core record. |
|---|---|
| Geography Structure Variable | Many data files that use geography as case identification will have a variable that will define the type of geographic level for a specific record (such as Country, county, city, etc.). This is a reference to that variable, which should provide a listing of each geographic type available and its type identification. Even without a full geographic structure description, this single variable can provide a great deal of information about the available geographies within the data file. |
| Spatial Object | Data are gathered and reported at a variety of object levels. This is one of the basic pieces of information geographers need to know in order to map data. Point data can be assigned or aggregated to lines, linear rings, and polygons, but polygons cannot be separated into their contained points. |
| Geographic Structure | Geographic structure provides detailed information on the types of geography available and how the various types relate to each other. Single hierarchies, layered hierarchies, and restricted coverage information are provided here. |
| Geographic Locations | When a file contains a small set of locations or in particular when it contains selected locations of a single type, it is often useful to have that information in the metadata. For example, does a file containing data on cities of 25,000 or more include the city of Black Duck, Minnesota?  Geographic |

| | locations may contain a more detailed bounding polygon description and/or pointers to external geographic shape or boundary files. |
|---|---|
| Summary Data Reference | If the geographic structure is defined, this element can be used to identify all levels that have summary data attached to them. For example, a file with the geographic structure of State/County/Tract may have data records for all levels or only for the lowest level (higher levels are created through aggregation). |
| Top Level Reference | This identifies the broadest area of coverage by text name and also optionally as a reference to the geographic level that describes it. |
| Lowest Level Reference | This identifies the smallest area of coverage by text name and also optionally as a reference to the geographic level that describes it. This is similar to the geographic unit as described in earlier versions of DDI. |

361
362  The Geographic Structure describes the levels of geography and their
363  relationships by defining a level and its parent or parents. Note that coverage
364  limitation information can be included at specific levels to further define the
365  coverage structure. This information will help external systems identify when only
366  a subset of locations are available in the data (such as only Counties with 100 or
367  more farms of $1,000 or more in yearly income) and the type and location of
368  geographic codes in the data.
369
```
370  <r:Geography isIdentifiable="true" id="GEO_0">
371          <r:Level>
372                  <r:Code>010</r:Code>
373                  <r:AurhorityOrganziationReference isReference="true" isExternal="false">
374                          <r:ID>USCB</r:ID>
375                  </r: AurhorityOrganziationReference>
376                  <r:Name>Country</r:Name>
377          </r:Level>
378  <r:Geography>
379  <r:Geography isIdentifiable="false" id="GEO_1">
380          <r:Level>
381                  <r:Code>040</r:Code>
382                  <r:AurhorityOrganziationReference isReference="true" isExternal="false">
383                          <r:ID>USCB</r:ID>
384                  </r: AurhorityOrganziationReference>
385                  <r:Name>State</r:Name>
386          </r:Level>
```

```
387          <r:ParentGeography isReference="true" isExhaustiveCoverage="true">
388                  <r:ID>GEO_0</r:ID>
389          </r:ParentGeography>
390   <r:Geography>
391   <r:Geography isIdentifiable="false" id="GEO_2">
392          <r:Level>
393                  <r:Code>050</r:Code>
394                  <r:AurhorityOrganziationReference isReference="true"  isExternal="false">
395                          <r:ID>USCB</r:ID>
396                  </r:AurhorityOrganziationReference>
397                  <r:Name>County</r:Name>
398                   <r:CoverageLimitation>Counties with at least 100 farms with $1,000 or in
399                  more yearlyincome</r:CoverageLimitation>
400          </r:Level>
401          <r:ParentGeography isReference="true isExhaustiveCoverage="true">
402                  <r:ID>GEO_1</r:ID>
403          </r:ParentGeography>
404   <r:Geography>
405
```

The complementary piece to this is Geographic Locations where explicit
locations can be listed. This is useful for files with limited geography or where
shape files will be attached to specific geographic locations. The description
starts with information on a specific type of geography which indicates if there is
any coverage limitation, and then identifies the variable containing the
geographic code for that area type and the authority reference for the coding
system. This is followed by a list of specific location codes. Note that specific
codes do not need to be included, so that you can provide all of the information
noted above without going into further detail.

Individual location listings provide a specific geographic code, a name, and the
geographic time (valid time period for the geography being used as this does not
always match that of the data), as well as the option of including a bounding
polygon (and excluding polygon) and/or reference to a shape file to describe the
area in as detailed a manner as desired.

```
422   <r:GeographicLocation isVersionable="true" id="GL_1">
423      <r:AurhorityOrganziationReference isReference="true" isExternal="false">
424              <r:ID>USCB</r:ID>
425          </r: AurhorityOrganziationReference>
426        <r:GeographicLevelRefereence isReference="true">
427        <r:ID>GEO_2</r:ID>
428      <r:GeographicLevelRefereence>
429      <r: Values>
430        <r:VariableReference>
431          <r:Reference><r:ID>STCNTY</r:ID></ r:Reference>
432        </r:VariableReference>
433        <r:GeographyValue>
434          <r:GeographyCode>
435              <r:Value>2700010</r:Value>
```

```
436         </r:GeographyCode>
437         <r:GeographyName>Aitkin [MN]</r:GeographyName>
438           <r:GeographicTime>
439         <r:StartDate>1860</r:StartDate>
440         <r:EndDate>9999</r:EndDate>
441           </r:GeographicTime>
442           <r:BoundingPolygon>
443             <r:ExternalURI>http://data.nhgis.org/MNcty1860.prj</r:ExternalURI>
444             <r:PolygonLinkCode>2700010_1860</r:PolygonLinkCode>
445             <r:GeographicTime>
446               <r:SimpleDate>1860</r:SimpleDate>
447             </r:GeographicTime>
448           </r:BoundingPolygon>
449       </r:GeographyValue>
450       <r:GeographyValue>
451         <r:GeographyCode>
452           <r:Value>2700310</r:Value>
453         </r:GeographyCode>
454         <r:GeographyName>Cook [MN]</r:GeographyName>
455         <r:GeographicTime>
456           <r:StartDate>1880</r:StartDate>
457           <r:EndDate>9999</r:EndDate>
458         </r:GeographicTime>
459       </r:GeographyValue>
460     </r:Values>
461   </r:GegraphicLocation>
```

# 2.0  Technical Structures

DDI 3.0 is dependent upon creating relationships by reference. While a hierarchy
may seem more intuitive to a user, a strict hierarchical approach limits the
availability of materials for reuse. Reuse facilitates the development of
documentation throughout the life cycle of the study, reduces the possibility for
human entry error, and provides a basic level of implicit comparability. In order to
provide interoperability for reference systems (allowing for exchange and reuse
of existing metadata) a consistent form of identification must be employed.

## *2.1  Identifiable Objects*

Not all elements are identifiable, but identification, when applicable, is required. If
we want DDI metadata to be shared and act as a transport structure to run
programs and processes, consistent use of identifiers is essential. Non-identified
elements generally require context and "live" within complex elements that are
identifiable and provide context for the non-identified element.

Elements that have identification only, that are not versionable, inherit their
version from their versionable parent, and their agency from their maintainable
parent agency. In other words, if the version or agency of the parent element
changes the identifiable element is considered to be part of this new version.

481  Elements that are identifiable only use the complex element r:IdentifiableType as
482  an extension base. This provides it with a consistent set of attributes that are
483  used to identify the element and to note any local changes for inherited contents.
484  The element may also be assigned a Name which can be repeated for language
485  and/or geographic alternatives. All identifiable elements must have an ID entered
486  in the attribute id. The full identification can also be expressed as a URL which
487  includes the complete path name. See Part I Appendix 1 for the full URL paths
488  for all identified objects in DDI 3.0. All identifiable elements within a DDI instance
489  can be located by the required fixed attribute isIdentifiable="true" and contain the
490  following element and attributes (attributes are identified with the prefix @ and all
491  start with lower case letters):

492

| Element / attribute | Description of use |
| --- | --- |
| @id | Required identifier for the element. This MUST be unique within the parent maintainable. |
| @urn | An optional urn for the element. Note that if there is conflict between the id and urn content, the urn takes precedence. |
| Name | An optional repeatable element which allows for a human-readable name to be attached to the parent element. It can be repeated for language and/or geographic alternatives. |
| @action | Action has a controlled vocabulary of "Add", "Update", or "Delete". It is used to identify local overrides to inherited content.<br>• *Add* – the element has a unique id and should be used in addition to the inherited elements<br>• *Update* – the element has the id of the inherited element which it updates (for example a local Name or label change)<br>• *Delete* – the element has the id of the inherited element which is NOT used in the local instance (for example a ProcessingEvent was not used). Note that if the identified element is complex, the entire contents of the complex element will be considered as deleted. |

493
494  Note that the attribute action is used only with inherited materials. Inheritance
495  occurs with grouping. These action statements provide local overrides for the
496  current inheritance, they themselves cannot be inherited. Note that if the element
497  that contains the change is not identifiable, its parent identifiable should be
498  entered in full including the changed information.

499  ## *2.2   Versionable Objects*
500  A subset of identifiable elements is also versionable. These are elements for
501  which changes in content are important to note. They use the extension base of
502  VersionableType and include all of the elements and attributes of
503  IdentifiableType. They are identified in a DDI instance by the fixed attribute

504  isVersionable="true". In addition Versionable elements allow one to indicate a
505  major and minor version, the date of the version, who changed it, and why it was
506  changed. If a change occurs at a lower level, it requires a version change of its
507  parent. Users need to understand if the change that was made will affect their
508  analysis of the data. To do this, they need to know why a change was made and
509  who made it. Versionable elements include the following elements and attributes
510  in addition to those listed under Identifiable.
511

| Element / attribute | Description of use |
| --- | --- |
| @version | If this is missing the default is assumed to be version 1.0. This attribute can contain only numbers and the separator "." But can extend to as many levels of specificity as needed by the maintenance agency. The first number to the left of the separator is the major version number. All subsequent numbers are considered minor versions. |
| @versionDate | This is the date that the specific version becomes active. It is a simple date containing a dateTime, date, YearMonth, or YearMonth. It should be as specific as possible. |
| VersionResponsibility | Do not use this to indicate the name of the Maintenance Agency as this is the only agency which can make changes in the content of the instance. VersionResponsibility was provided to allow for the identification of the person or suborganization within the maintenance agency that made the change. This may be important to internal management. |
| VersionRationale | This provides the reason for the change. The correction of a typographical error may have different ramifications for the end user than the replacement of erroneous content. |

512
513
514  Note that versioning is only required for published materials. Maintaining version
515  identification during the production process will depend on the needs of the
516  organization producing the document.

517  ## 2.3  Maintainable Objects
518  There are a number of complex elements that represent major blocks of objects
519  that can be maintained outside of a DDI Instance (published as separate
520  entities). All maintainables are published within a DDIInstance under one of the
521  following packaging structures: StudyUnit, Group, or ResoursePackage. Note
522  that a PhysicalInstance can only be published within a StudyUnit or Group. Major
523  modules can be maintained as can all complex elements whose name ends in
524  "Scheme", plus occasional additional complex elements. Note that
525  ResourcePackage and Group are the two possible top level elements of the

526  module scheme Group. Most module schemes contain a single top level
527  containing element which carries the name of the module. Version changes can
528  only be made by the maintaining agency (or at their specific request) so that if a
529  non-maintaining agency makes a change.
530

| Modules | Scheme | Additional |
|---|---|---|
| DDIInstance | QuestionScheme | Instrument |
| Archive | CategoryScheme | ResourcePackage |
| StudyUnit | CodingScheme | Group |
| DataCollection | VariableScheme | |
| LogicalProduct | ConceptScheme | |
| PhysicalDataProduct | UniverseScheme | |
| PhysicalInstance | OrganizationScheme | |
| Comparative | GeographicStructureScheme | |
| ConceptualComponent | GeographicLocationScheme | |
| DDIProfile | NCubeScheme | |
| | PhysicalStructureScheme | |
| | RecordLayoutScheme | |
| | ControlConstructScheme | |
| | InterviewerInstructionScheme | |

531
532  Maintainable objects can be used to assemble an instance by reference or to
533  share commonly used sets of information. Elements that are maintainable are
534  extensions of the complex element r:MaintainablType. They include all the
535  elements and attributes of VersionableType and add the attribute agency. The
536  content of Agency is the Name token of maintenance agency. This unique token
537  must be declared in the document as a full organization or individual description
538  or as a reference to an external registry of organizations. Note that the id of all
539  maintainable objects MUST be unique within the maintaining agency.
540

## 541  *2.4  Constructing an Identification*

542  The complex elements IdentifiableType, VersionableType, and MantainableID
543  are used to identify described objects for the purposes of internal and/or external
544  referencing. All identification allows for two ways of expressing an entity's
545  identification. The first is through the use of a URN. The URN provides the
546  element type, the maintaining agency, version, and element ID. While URN is
547  recommended, the alternative is using a series of elements that provide the
548  same information as found in the URN. If both URN and a combination of
549  elements are used, the URN takes precedence. The DDI uses a specifically
550  structured URN [see URN part 2.5]. Note that VersionableTypes inherit their
551  IdentifyingAgency from their parent Maintainable object, while the
552  IdentifiableType inherits both its versioning information and IdentifyingAgency
553  from its parent Versionable object. Any object that is published (a Group,
554  StudyUnit, or any Maintainable Object published for inclusion by reference, must

555 be published within a DDIInstance. The ID of this instance MUST be unique
556 within the maintaining agency.
557
558 Identification contains a URN and/or the sequence of elements shown below. In
559 addition Identification has an optional repeatable element Name. This is a
560 human-readable name given the entity being identified. This may be a mnemonic
561 that is commonly related to the element. It may be repeated to provide language
562 and/or geographic alternatives.
563

| ELEMENT/ATTRIBUTE | TYPE | Description |
|---|---|---|
| @id | NCName | ID assigned by an agency. This must be a unique identifier within the maintained object. It must start with a character and can contain any character or number plus any of the following non-alphanumeric characters: "*", "@", "_", "$", or "-" |
| @agency *(MaintainableType only)* | NCName | The agency maintaining the identification. This is optional if inherited from a parent object. The published instance must contain an entry in an OrganizationScheme either inline or by reference. The content should be the ID of the organization/individual. |
| @version *(VersionableType and MaintainableType only)* | string | Version number, expressed as a two-part numeric string composed of two positive integers separated by a period. The first number indicates a major version, the second a minor one: 1.0. Optionally, a third integer may indicate sub-version: 1.0.2. |
| @versionDate *(VersionableType and MaintainableType only)* | BaseDate | Date the version took effect expressed as dateTime, date, YearMonth, or Year |
| VersionResponsibility *(VersionableType and MaintainableType only)* | string | Reference to the person and/or organization responsible for the version change. This is primarily intended for internal use and can be as detailed as the organization requires. |
| VersionRationale *(VersionableType and MaintainableType only)* | string | Textual description of the rationale/purpose for a version change. This should be informative to users so that they can determine if the change has potential impact on their analysis. |

564

565

## *2.5   URN Structure*

567 A DDI URN has a specific structure that must be followed for uses of the URN in
568 identification and reference to DDI objects.
569 **urn="urn:ddi:3_0:<Maintainable Object Class.Object Class>=<Agency ID>:<ID**
570 **of maintained object>[<Major Version>.<Minor Version>].<ID of contained**
571 **object>"**
572
573 An example of the URN structure for a versionable object is broken down below:
574
575 **urn:ddi:3_0:VariableScheme.Variable=ICPSR:VScheme_4[1.0]AGE_3[1.0]**
576

577 • The delimiters are as follows
578   o Top level field separator   **:**
579   o Hierarchical separator      **.**
580   o Object class to identification separator   **=**
581   o Version separator   **[ ]**
582 • All sections are required.
583 • All published elements are either maintainable or contained within a
584   maintainable object. The maintainable object must be identified first, as
585   the ID is unique within the maintainable. The object type identification will
586   always contain a maintainable object plus the optional specific object if it is
587   either Versionable or Identifiable
588 • The full path of ids with version where applicable should be provided (Part
589   I Appendix 1 contains an alphabetical listing of identifiable and versionable
590   elements in alphabetical order with their complete nesting string up to
591   parent maintainable)
592
593 Breakdown of the example URN:
594

| Identify that this is a URN | urn |
|---|---|
| Top level field separator | **:** |
| It is a DDI URN | ddi |
| Top level field separator | **:** |
| It is DDI version | 3_0 |
| Top level field separator | **:** |
| The maintainable object type | VariableScheme |
| Hierarchical separator | **.** |
| The object | Variable |
| Object class to identification separator | **=** |
| The maintenance agency is the DDI Alliance | ICPSR |
| | |
| The name of the VariableScheme is | VScheme_4 |

| The version number of this element is encased in version separator | [1.0] |
|---|---|
| Hierarchical separator | . |
| With the ID of object | AGE_3 |
| The version number of this element is encased in version separator | [1.0] |

595

596 Note that the Variable must have the same maintenance agency as the
597 maintainable it resides in, and therefore this is not repeated.

598

599 If the element being referenced was a maintainable itself, then the URN would
600 end after the first version element.

601 *Examples*

602 **URN of a maintained object**

603 To identify of a variable scheme in DDI 3.0 via a URN would be as follows:
604 **urn="urn:ddi:3_0:VariableScheme=ICPSR:V_GENDER_SCHEME[1.0]"**.

605 **URN of an versionable object**

606 All versionable objects are contained within maintainable objects. To identify of a
607 variable in DDI 3.0 via a URN would be as follows:
608 **urn="urn:ddi:3_0:VariableScheme.Variable=ICPSR:V_GENDER_SCHEME[1.0].Male**
609 **[1.0]"**

610 **URN of an identifiable object**

611 An identifiable object may be a direct child of a maintainable object or be
612 contained by a versionable object within a maintainable object. The full path
613 should be provided to facilitate locating the item when referenced.

614

615         &lt;DataCollection isMaintainable="true" id="DC_5698" version="2.4"&gt;
616             &lt;Methodology isVersionable="true" id="Meth_Type_1" version="1.0"&gt;
617                 &lt;TimeMethod isIdentifiable="true" id="TM_1"&gt;

618

619 To identify the identifiable object in the above hierarchy in DDI 3.0 via a URN
620 would be as follows:
621 **urn="urn:ddi:3_0:DataCollection.TimeMethod=ICPSR:DC_5698[2.4].Meth_Type_1[**
622 **1.0].TM_1"**

623 **URN of an object that nests within its own object type**

624 An example of this is an Individual who belongs to an Organization that is nested
625 in another Organization. In this case each object type would be listed in order
626 and the IDs of the full path would be provided in the URN.

627

628         &lt;OrganizationScheme isMaintainable="true" id="OS_1" verson="1.0"&gt;
629             &lt;Organization isVersionable="true" id="UMICH"&gt;

630          &lt;Organization isVersionable=&quot;true&quot; id=&quot;ICPSR&quot;&gt;
631                &lt;Individual isVersionable=&quot;true&quot; id=&quot;J_Doe&quot;&gt;
632
633   **urn=&quot;urn:ddi:3_0:OrganizationScheme.Individual=ICPSR:OS_1[1.0].UMICH[1.0].IC**
634   **PSR[1.0].J_Doe[1.0]&quot;**
635

## *2.6   Referencing*

637   DDI 3.0 contains four types of references:
638

| | |
|---|---|
| Reference | The basic structure used to reference DDI objects |
| SchemeReference | A special extension to Reference for referencing DDI schemes |
| OtherMaterial | The basic structure for referencing external non-DDI objects |
| Image | Reference to an external image file currently used only in the OrganizationScheme |

639
640
641   References to DDI objects are all based on the ReferenceType as described in
642   reusable. Elements within a DDI instance which reference DDI objects will all be
643   identified with the fixed attribute isRefrence=&quot;true&quot;. Note that when referencing an
644   object internal to the DDI instance you may simply provide the ID as long as it is
645   unique within the instance. However, if this instance is later parsed into its
646   constituent parts, the parser will need to supply the missing components to
647   uniquely resolve the reference. If the ID of an object is only unique within its
648   maintainable parent (module or scheme), the identification of that maintainable
649   parent must be included to provide a unique identification. Once again, if both a
650   URN and ID are provided, the URN takes precedence if they have conflicting
651   content. Reference includes the following structure.
652

| Element/Attribute | Description |
|---|---|
| Module | This is a complete reference structure for the parent module. This must be used in cases where there have been local modifications. [optional] |
| Scheme | This is a complete reference structure for the parent scheme.  [optional] |
| URN | A DDI structured URN may be used in addition to or instead of an ID sequence. |
| ID | ID of the element being referenced. This is the minimal required content for a reference in terms of content and length. |
| IdentifyingAgency | The NCName of the maintenance agency used as part of the ID sequence [option] |

| Version | The version number of a referenced versionable or maintainable object used as part of the ID sequence [optional] |
|---|---|
| @isExternal | The default setting is "false", indicating that the element will be found inline within the DDI instance. If this value is set to "true" (reference to an external DDI object) the URI must be provided. |
| @URI | This is the URI for the external location when isExternal is set to "true" |
| @isReference | Fixed Boolean value of "true" |
| @lateBound | A Boolean attribute with the default setting of "false" indicating that the reference is to a specific version of the object. When changed to "true" the reference will retrieve the most recent version available. Note that while identifiable objects do not have a version themselves, the object will be obtained from the most recent version of its parent versionable or maintainable object. |

653
654    The use of late bound and early bound references provides flexibility in the
655    material obtained by the referernce. Early bound (@lateBound="false" means
656    that you are referencing a specific version and that none other will do. This
657    allows you to replicate the metadata as it was used in a previous analysis. Late
658    bound allows you to request the most recent version of the referenced element.
659    For example, you may always wish to include the most recent version of an
660    Archive module to receive the most current updates of organizational information,
661    life-cycle events, related publications, and access information. Or you may want
662    to be sure you have the most current corrections and updates to the information
663    in a logical product.
664
665    SchemeReference extends Reference and provides an addition element set that
666    allows you to exclude one or more identifiable elements from the scheme being
667    included by reference. It allows a user to include say a VariableScheme but
668    exclude one or more variables within that scheme. Most schemes will already
669    allow for selective inclusion of component items from other schemes, but this
670    particular reference simplifies the process when the uses wishes to include the
671    majority of the scheme, excluding only a few items. The additional element
672    Exclude contains an ID and Version element pair where the version is optional.
673    Since the parent maintainable scheme has already been identified, this is all that
674    is needed to identify specific items for exclusion.
675
676    OtherMaterial as a reusable is available in all major modules. References to
677    related materials, physical objects, etc. are listed here and can be attached to
678    any identifiable object. This material is not intended to be processed by DDI but
679    is supplied to inform the end-user about non-DDI materials related to the

680    development, collection, analysis, or other materials related to the DDI Instance.
681    OtherMaterial contains the following:
682

| Element/Attribute | Description |
|---|---|
| Citation | Full citation element for the external object |
| ExternalURLReference | Location of an external electronic object via a URL (for example http://icpsr.umich.edu/ ) |
| ExternalURNReference | The namespace of the external object expressed as a URN (for example urn:isbn:0-395-36341-1 the unique International Standard Book Number) |
| Relationship | This contains both a reference to an identified DDI object and RelationshipDescription (optional/repeatable). Relationship can be repeated to link to multiple DDI Objects. |
| MIMEType | A standard internet MIME type for use by processing applications |
| @type | A required type code for type of OtherMaterial. |

683
684    External objects that are not DDI require the more complex structure of
685    OtherMaterial. The use of non-DDI references is quite limited outside of the
686    standard section for OtherMaterial listings within each of the major modules.
687    Currently use is restricted to ExternalInterviewerInstruction, ExternalAid, and
688    r:OtherMaterial in Generation. ExternalInterviewerInstruction provides the ability
689    to reference non-DDI structured interviewer instructions from with a
690    ControlConstruct for a questionnaire/instrument. It consists of OtherMaterial and
691    an attribute displayText which indicates whether the information should be
692    displayed or simply made accessible (via reference or specific request for
693    display). ExternalAid is available in QuestionItem, MultipleQuestionItem, and
694    ControlConstruct. Its purpose is to provide access to an image, sound, media, or
695    physical object used in the question. For example "Watch the 30 second ad and
696    indicate if you ….." where ExternalAid would be the link to the advertisement.
697    ExternalAid is type="r:OtherMaterial". Generation is a structure to describe how a
698    specific Category within a CategoryScheme was generated and includes a child
699    element r:OtherMaterial. An example of how this might be used is a reference to
700    a printed table of Poverty breakpoints for the United States definition of poverty.
701    The Generation would describe how the breakpoints were generated and then
702    provide a link to the resultant table.

### 2.7   *Date*

704 The DDI provides a structure for a Date element which allows a choice between
705 single, simple dates or date ranges. If a date element contains a range, Cycle
706 may be used to indicate occurrence of this range within a series of ranges. An
707 integer is used to identify the cycle. Dates are required to be expressed as ISO
708 860-formatted dates for all fields. The optional attribute calendar on any parent
709 date item, such as PublicationDate, allows you to designate a non-standard
710 calendar type. Historically-formatted dates may be included in addition for
711 archival or other purposes.
712

### 2.7.1  Simple Date

714 <r:PublicationDate calendar="Georgian">
715 <r:SimpleDate>2007</r:SimpleDate>
716 </r:PublicationDate>
717
718 This is a single point in time and conforms to any of the following ISO 8601
719 standard structures.
720

| dateTime | YYYY-MM-DDThh:mm:ss | 1982-01-05T23:05:15 |
|---|---|---|
| date | YYYY-MM-DD | 1982-01-05 |
| gYearMonth | YYYY-MM | 1982-01 |
| gYear | YYYY | 1982 |
| duration | PnYnMnDTnHnMnS | P26Y2M22DT11H5M20S |

721
722 Note that the "T" in dateTime is literal, denoting the beginning of the Time
723 section, and that "ss" can contain decimals. Optionally, dateTime can be
724 extended by a timezone offset of "Z" to represent Zulu time or GMT. For
725 example, Eastern Standard Time is Z-4.
726
727 Note that the "P" in duration is literal and indicates that this is a Period of
728 duration. The other upper case letters are also required with the preceding
729 number providing the number of years (nY), months (nM), etc. A period may be
730 of negative duration, for example a period of minus 10 days (-P10D), by
731 preceding the "P" with a negative sign.
732

### 2.7.2  Date Range

734 A range is expressed as a StartDate and EndDate each expressed in the same
735 format as a simple date. The dates are assumed to be inclusive. The position of
736 this range within a series of ranges is expressed as an integer in Cycle. For non-
737 standard calendars an attribute calendar allows specification of the calendar
738 used.
739
740 <r:Date>
741 <r:StartDate>2006-04-01</r:StartDate>

742    <r:EndDate>2007-03-31</r:EndDate>
743    </r:Date>

## 2.7.3 Historical Dates (expressed in formats other than ISO 8601)

745    All dates can optionally be expressed in other historical structures with an
746    attribute to describe the structure being used. This is simply a string containing
747    the historical date and an attribute historicalDateFormat used to specify the non-
748    ISO date format. For example:
749
750    <HistoricalDate historicalDateFormat="Month DD, YYYY">January 5, 1982</HistoricalDate>
751
752    Historical date information parallels the simple date, start date and end date
753    structures of the standard DateType.
754

## *2.8   String Types*

756    String or text entries have a number of formats to support language differences,
757    the need for structured text, and constraints on content.
758
759

| Sting Type | Features |
|---|---|
| NCName | Must start with a letter and can contain alphanumeric \| "_" \| "." |
| String | Any character string (will be read as the literal string) |
| InternationalString | A string with an xml:lang attribute to denote language and boolean attributes translated (default false) and translatable (default true) |
| StructuredString | In addition to features of InternationalString allows for XHTML structure tags in the content |
| IdentifiedStructuredString | Combines features of IdentifableType and StructuredString |
| DynamicText | Structures the behavior of dynamic or static text within a question by allowing a text line to be broken into segments describing both static (literal text) and dynamic (conditional text ) . (See section 5.1 for details of use) |

760
761    The following grid shows which features are available for each type other than
762    NCName. Many of the forms without ID are parts of complex elements that are
763    identifiable.
764

| | string | ID | xml:lang | translated | translatable | XHTML |
|---|---|---|---|---|---|---|
| String | X | | | | | |
| InternationalString | X | | X | X | X | |
| StructuredString | X | | X | X | X | X |
| IdentifiedStructuredString | X | X | X | X | X | X |

765

## *2.9   DDI Profiles*

767   Many DDI users first determine which elements they will use, or not use, and
768   determine any constraints on how they are used for their specific local needs.
769   This is particularly common in organizations who wish to make sure elements are
770   used in a consistent manner, or need information expressed in a specific way to
771   support their systems. Common examples of this are the DDI Core Elements, the
772   CESSDA list of required elements for inclusion in their search system, or NHGIS
773   rules for special use elements. Each of these examples imposes additional
774   required elements, informs the user what items are expected in submitted
775   documents, restricts what is handled by a system, or specifies accepted usage of
776   an element within a system.

777

778   DDI profiles allow this type of specification to be defined in a consistent way that
779   can be published and used for validation. This profile describes which DDI
780   elements and attributes are used in a particular use profile or supported by a
781   DDI-conformant application. It uses XPath expressions to identify the used or
782   unused fields in terms of the full possible DDI instance. Its construction is quite
783   simple, but it allows an organization to require certain elements for a DDI
784   instance to be included in a collection (example, CESSDA), to inform contributors
785   of system limitations (example, does not handle NCubes), or to inform an
786   organization on the proper use of elements for internal compliance.

787

| Object | Purpose |
|---|---|
| Identification | Identifies the DDI Profile. |
| XPathVersion | Provides the version of XPath used. Currently values are either 1.0 or 2.0. |
| DDINamespace | This is the DDI version, currently 3.0 |
| Used | A repeatable complex element that describes a DDI element used by the profile. This complex element contains an XPath, which points to the element or attribute that is being used. All sub-elements of a used element are assumed to be supported unless explicitly addressed by the profile. The number of supported repetitions may be included in the XPath. An attribute "required" with a default value of "false" allows for requiring elements that are optional in the DDI schema. Note that if a field is required by the DDI Schema, it cannot be made optional. |
| NotUsed | A repeatable complex element that describes a DDI element NOT used by the profile by providing its XPath. A required DDI Element cannot be disallowed. |

788

789   Used allows the identification of an AlternateName for the object, a Description
790   and Instructions for use. The attributes include required for changing an optional
791   object to a required object, path, defaultValue to designate a default value if the

792    object is missing, and a fixedValue indicating that the defaultValue noted in the
793    earlier attribute is actually a fixed value.

## 794    3.0   Capturing the Background Information

### 795    *3.1   Study Unit*

796    The first stage in capturing information for a study is always within the study unit.
797    The study unit captures information related to the idea behind the study, the early
798    stages of proposal or project development, and provides it with an identification
799    and basic information on scope and coverage. Remember that this is preliminary
800    content and can be altered as the study progresses through the life cycle. DDI
801    3.0 was designed to be used for capturing information across the life cycle
802    process of the study (as well as after the fact) and can be used to hold drafts,
803    ideas, and notes that will be formally incorporated at a later date. Remember that
804    once "published," any changes must be versioned.
805
806    Start a DDI xml instance using studyunit.xsd. In essence this section captures
807    content from the inception of the study, from the idea. Include references to all
808    other schemas you may be using, but at minimum include the following xmls
809    identifiers:
810

| .xsd file | Reason for inclusion |
|---|---|
| studyunit | This is the base source xsd. |
| reusable | All schemas use elements from reusable.xsd. |
| conceptualcomponent | Concept, Universe, and Geography definitions are some of the earliest sets of information created. |
| archive | Captures information on your organization and process and identifies the maintaining agency. |
| datacollection | Allows you to capture information on planned methodology such as sample design. |

811
812    Complete the Identification and use all the fields (id, agency, and version) as
813    these can be inherited by elements lower in the structure.
814
815    Under Citation provide at least a working Title and Creator.
816
817    Use Abstract to capture an outline of the study. Remember that this can be
818    edited later to relate further refinements. The null hypothesis or other description
819    of why the study is being conducted (intended use) should be placed in Purpose.
820
821    Complete the upper level or levels of Universe in conceptualcomponents and
822    provide a link from the study unit to the uppermost level of the universe (see
823    Universe, Part 3.5). This is also a good point to begin collecting and structuring
824    your concepts (see Concepts, Part 3.6). These two sections are needed by later
825    modules for reference purposes. While they can be edited and expanded at any
826    point, completing them at this point assists you in clarifying both the population

827 being studied and the concepts covered by the study. The
828 GeographicStructureScheme and GeographicLocationScheme can be started
829 now if known. If external ResourcePackages are available which contain a
830 standard geographic structure and location coding system, it is worthwhile to
831 note it now by referencing those schemes. This information will inform how
832 questions will need to capture geographic information for the study. Include your
833 unit of analysis in AnalysisUnit.
834
835 Coverage at this level is the overall coverage of the study. You must define the
836 coverage at this level, as it is inherited or constrained by all lower levels. Lower
837 levels can constrain the coverage by limiting the topical, temporal, or spatial
838 coverage of a specific module, but they cannot expand the coverage. Coverage
839 is made up of the following sections and at least the minimal level information
840 should be provided.
841

| COVERAGE | PRIMARY TAGS | PURPOSE OF INFORMATION |
| --- | --- | --- |
| Topical | Subject, Keyword | Provides structured subject headings and unstructured keywords |
| Temporal | ReferenceDate | Reference dates provide information on dates covered by the data that are outside of the collection period, such as residence 5 years ago. |
| Spatial | Description, SpatialObject, TopLevelReference, LowestLevelReference | While this may be expanded later in the life-cycle of the study, these are the minimal level recommended tags. Description is the equivalent of the Dublin Core geographic coverage. Spatial Object at this point is the most discrete level at which data will be captured (point, line, polygon, or linear ring). Top Level and Lowest Level are the broadest area of spatial coverage (for example, Europe) of the study, and the most discrete identifiable level (for example, a NUTS 3 or a housing unit). These three pieces of information are used by geographers to determine the usability of the contents for various GIS uses as well as providing basic geographic |

| | | structural information to the data user. |
| --- | --- | --- |

842
843
844  Begin building the OtherMaterial section of the studyunit at this point. Use it to
845  capture citations on materials used as references in developing the study
846  proposal. Provide a citation for any proposals submitted to funding agencies and
847  any other related material. Some of this may be incorporated inline at some
848  point, but you will still want to retain the citation to the separate document.
849
850  "Note" can be used to capture information you may want to include in other areas
851  at a later date, but want to hold in a temporary area of your instance. Select an
852  appropriate type attribute so that at a later date these will be easy to locate,
853  attach to appropriate elements, or transfer information to another location and
854  delete.

855  ## 3.2  Concepts

856  Concepts for variables and questions are described in a concept scheme and
857  referenced by the variable or question. A question can reference multiple
858  concepts; however, a variable must be linked to a single concept. Variables that
859  capture a number of concepts require the creation of a composite concept that
860  captures both concepts under a single concept description. This tends to occur
861  only in some legacy aggregate data files which nest a number of concepts in a
862  non-regular hierarchy. For example "Household Type by Presence and Age of
863  Related Children" would require a concept that includes both a concept for the
864  household type and one or more to capture presence and age of related children.
865  Alternatively, the variable can be split into its component parts (separate
866  concepts) and reassembled within an NCube that defines those sections without
867  data (see NCubes Section 7.5 for details).
868
869  In terms of a complete variable description, the concept provides the property of
870  the object in a complete ISO/IEC 11179 description. DDI uses the variable to link
871  the object (universe) with the property (concept) and the representation. A
872  ConceptScheme with just a collection of concepts requires the use of the variable
873  to tie these pieces together. Concept scheme can also contain data element
874  concepts which provide these links internally so that a complete ISO/IEC 11179
875  description can be held as a resource package outside of the context of a dataset
876  and its variables. A concept is versionable and contains a label, description, and
877  information on similar concepts (reference to the similar concept and a
878  declaration of the difference between the two). A concept may be identified as a
879  "characteristic" through the use of a Boolean attribute. This attribute must be set
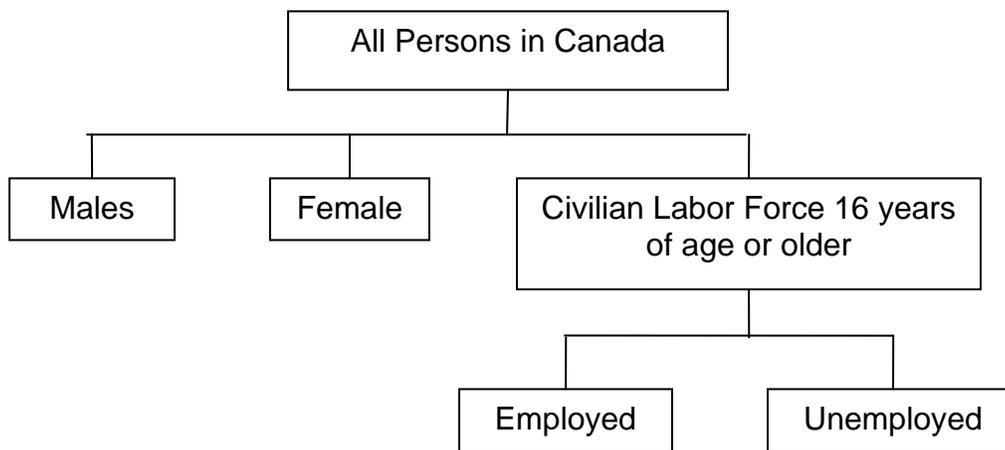880  to "true" if it is referenced by a data element concept.
881
882  A data element concept is a concept in the ISO/IEC 11179 sense. It contains a
883  label and description and can identify similar concepts like a general concept.
884  However, it references another concept which has its "isCharacteristic" attribute

885   set to "true" as well as referencing a universe. In this manner it is stating that the
886   data element concept defines a concept that is characteristic of a specific
887   universe.

888   ### *3.3   Universe*

889   A universe scheme contains a structured listing of all the universes within a
890   study. The top level should reflect the complete universe of the study, for
891   example, "All persons and housing units in Canada". Lower level subsets of this
892   universe should be organized hierarchically. Siblings are not assumed to be
893   mutually exclusive, but a child must be a subset of its parent. For example,
894

```
                        ┌──────────────────────────┐
                        │   All Persons in Canada  │
                        └──────────────────────────┘
              ┌────────────────┬──────────────────────────┐
         ┌─────────┐      ┌──────────┐      ┌──────────────────────────┐
         │  Males  │      │  Female  │      │ Civilian Labor Force 16  │
         └─────────┘      └──────────┘      │   years of age or older  │
                                            └──────────────────────────┘
                                            ┌─────────────┬─────────────┐
                                       ┌──────────┐   ┌────────────┐
                                       │ Employed │   │ Unemployed │
                                       └──────────┘   └────────────┘
```

895
896
897   The universe as described in the scheme is referenced by the question construct
898   (use of a question within an instrument), variables, and NCubes. Describing
899   these in a single location provides information on hierarchical relationships, and
900   clear comparability of the universe of two questions or variables that reference
901   the same universe. In addition, a universe can be described in both human-
902   readable and machine-readable formats that provide explicit information on the
903   definition of the universe. Note that variables allow for multiple universe links.
904   When more than one universe is referenced, the universe is defined as the
905   common area. For example a Variable that references both the universe
906   "Female" and the universe "Employed" would be defined as Females in Canada
907   who have a status of Employed in the Civilian Labor Force (16 years of age or
908   older).

### 909    3.4    Methodology

910    If known at this point, include information on Methodology found in the
911    datacollection schema. Methodology has an extension base of VersionableType.
912    Given the breadth of the content, careful use of VersionRationale can help the
913    user to quickly identify material sections that were changed or added in the new
914    version. Methodology includes DataCollectionMethodology, TimeMethod,
915    SamplingProcedure, and DeviationFromSampleDesign. In addition, you can note
916    any datacollection software you plan to use in Software. Note that all of these
917    sections are repeatable and that you should note separate methods in separate
918    repetitions. You may need to reference a specific method later in the study and
919    you can only do this if you have listed and identified each as a separate
920    description. With the exception of Software, these are all
921    IdentifiedStructuredStrings to allow for clear presentation of the material.

### 922    3.5    Collection Event

923    Each collection event should have its own entry. A collection event is normally
924    described as the event (covering one or more days) that results in a raw data set
925    which may or may not be linked to other raw datasets over time (through
926    repetition) or from other sources for example a Parent survey linked to a Child
927    survey. Each CollectionEvent provides a reference to who collected the data
928    (DataCollectorOrganizationReference), a description of the DataSource,
929    DataCollectionDate, DataCollectionFrequency, ModeOfCollection,
930    CollectionSituation, and ActionToMinimizeLosses. While all of these except
931    DataCollectionDate are repeatable within CollectionEvent, if the details of
932    multiple CollectionEvents are listed within a single CollectionEvent section, the
933    relationships between specific elements of information such as DataSource and
934    ModeOfCollection can be lost.

935

| Element | Description of content |
|---|---|
| DataCollectorOrganizationReference | References to the organizations or individuals responsible for collecting data. References are to the descriptions contained in an Organization Scheme. |
| DataSource | Provides a SourceDescription (structured strin), a SourceType (may be from a controlled vocabulary), the Origin which is a citation or URI to an external non-DDI description of the data, and Characteristic such as the level of documentation available for the source data or other factors like water damage that could affect the quality of the data source. |
| DataCollectionDate | Standard date type allowing a single date, range, or duration. |

| DataCollectionFrequency | This is the intended frequency of data collection, where the current collection event lies, or is meant to lie within a series. It consists of both a standard date type and the ability to add a code or string such as monthly, yearly etc. |
|---|---|
| ModeOfCollection | IdentifiedStructuredString |
| CollectionSituation | IdentifiedStructuredString |
| ActionToMinimizeLosses | IdentifiedStructuredString |

936

# 937 **4.0 Archives, Organizations, and Life Cycle Events**

938 The Archive in the context of DDI 3.0 is the individual or organization responsible
939 for the DDI instance at a given point in time. At the onset of a study this may be
940 the principal investigator or their organization. The Archive module contains
941 information that is both persistent (travels with the metadata throughout its
942 lifetime) or is specific to the archive itself (information that has no relevance once
943 it is distributed to another archive). At the onset of a study, the primary pieces of
944 information entered in the Archive module include the identification of
945 organizations, a reference to the current archives as listed in organization, and
946 the addition of specific life cycle events as they occur. The contents of the major
947 sections are provided below for context.
948

## 949 *4.1 Archive Specific Information*

950 The first item is a reference to the current archive as listed in Organization.
951 Following this is a listing of the individual items that make up the instance as
952 viewed by the archive. This can include associated data files, referenced
953 schemes, etc. Each item in the list can have the following individual distinct
954 pieces of information. Alternatively, multiple items can be treated as a Collection
955 of Items or subordinate Collections. Flexibility was provided to support the
956 common uses by different archives at different stages of the life cycle.
957 DefaultAccess and FundingInformation for the specific Archive can be listed here
958 and will apply to the contents of the items or collections listed in this section. For
959 example this could include access being restricted to Faculty, Staff, and Students
960 of an academic institution. The contents of both Item and Collection are listed
961 below. Items that are unique to each are in italics to help the user determine the
962 most appropriate usage for their particular archive.
963

| Item content | Use |
|---|---|
| LocationInArchive | The location of an item within an archive, for example, if the archive needed to differentiate between multiple systems or different physical location options within the archive. For example a publicly accessible site or an off-line storage site. |

| Call Number | The specific identifier within this archive. |
|---|---|
| URI | The URN or URI of the item. |
| *Format* | *The format of the data or metadata holding.* |
| *Media* | *The current media of the holding.* |
| Study Class | Allows for an archive specific designation of a study class. |
| Access | Access information specific to this item |
| Original Archive Organization Reference | A reference to the source archive of the documented material. This is repeatable so that a chain of ownership can be recorded. |
| Availability Status | A human-readable description of availability. |
| Data File Quantity | Number of data files in the documented holding. |
| Collection Completion | A statement regarding the completeness of the documented holding. |
| Item | A subordinate Item |

964

| Collection content | Use |
|---|---|
| LocationInArchive | The location of an item within an archive, for example, if the archive needed to differentiate between multiple systems or different physical location options within the archive. For example a publicly accessible site or an off-line storage site. |
| Call Number | The specific identifier within this archive. |
| URI | The URN or URI of the item. |
| *ItemQuantity* | *The number of items in the collection.* |
| Study Class | Allows for an archive specific designation of a study class. |
| *DefaultAccess* | *Access information specific to the contents of this collection.* |
| Original Archive Organization Reference | A reference to the source archive of the documented material. This is repeatable so that a chain of ownership can be recorded. |
| Availability Status | A human-readable description of availability. |
| Data File Quantity | Number of data files in the documented holding. |
| Collection Completion | A statement regarding the completeness of the documented holding. |
| Item | A subordinate Item |
| *Collection* | *A subordinate Collection* |

965
966

967    Note that access and access restrictions are provided at a variety of levels. This
968    includes general access restrictions for the archive as well as specific access
969    restrictions to all or parts of the data. Conditions for use, disclaimers, citation and
970    deposit requirements and all related access forms are located here.
971    DefaultAccess under ArchiveSpecific would apply to the contents of the specified

972 archive. DefaultAccess at the Collection level applies to all the items within the
973 collection whereas item specific information is found in Item/Access.
974

| Element | Use |
|---|---|
| Confidentiality Statement | A structured statement |
| Access Permission | A form type including a statement of the access permission required, the form number and location, and an indication that the form is or is not required. |
| Restrictions | A structured statement |
| Citation Requirement | A structured statement |
| Deposit Requirement | A structured statement |
| Access Conditions | A structured statement |
| Disclaimer | A structured statement |
| Access Restriction Date | A standard date type allowing for a specified date or range. |
| Contact Organization Reference | A reference to and individual or office within an organization or to the organization itself. |

975
976 Finally, there is funding information. In this location it is limited to information on
977 funding sources specific to the archive, for example, if an archive received
978 funding to collect, process, and provide access to a specific set of studies.

### 979 *4.2 Organization*

980 The OrganizationScheme is comprised of two basic structures: organization and
981 individual. Each basic structure is made up of a description of the organization or
982 individual and its components, the role of the organization/individual as related to
983 the life cycle of the data, and relationship to other organizations or individuals.
984
985 The organization description provides each organization with an ID that can be
986 referenced by other locations in the DDI instance. The name and abbreviation (or
987 nickname) of the organization, along with contact information and member
988 individuals should be provided here. An organization can belong to a group of
989 Organizations directly as a child of the group or through inclusion by reference.
990 Organizations can be related to other organizations in any way that can be
991 described by the user of the DDI. Individuals can belong to organizations (as well
992 as have organizations attached to the individual's record internally). They carry
993 their own contact information as well as information on language capabilities.
994
995 Designed to provide a flexible structure to describe organizations and individuals
996 and their relationships to each other, the organization module captures
997 information on organizations and individuals related to the study in a uniform way
998 and allows the DDI document to reference these entries in relation to their
999 specific roles within the life cycle of the study. The advantage, in addition to
1000 reducing redundancy, is the ability to capture the relationships between and
1001 within organizations, thereby clarifying the relationships of individuals and
1002 organizations related to the study.

1003
1004    Minimal entry for beginning a study would be the identification of the organization
1005    proposing to do the study:
1006
1007    <a:OrganizationScheme isMaintainable="true" id="MAINORG" urn="
1008        urn:ddi:3_0:OrganizationScheme=MPC.MAINORG[1.0]">
1009            <a:Organization isVersionable="true" id="MPC" urn="
1010            urn:ddi:3_0:OrganizationScheme.Organization=MPC:MAINORG[1.0].MPC[1.0] >
1011                    <a:OrganizationName>Minnesota Population Center</a:OrganizationName>
1012                    <a:Nickname>MPC</a:Nickname>
1013                    <a:URL>pop.umn.edu</a:URL>
1014            </a:Organization>
1015    </a:OrganizationScheme>
1016
1017    Any organization that is producing or managing a number of DDI documents
1018    should consider maintaining a separate OrganizationScheme within a published
1019    ResourcePackage. This makes tracking and updating organization changes
1020    easier and provides a clear historical trail. Note that Organization and Individual
1021    have to places to enter a number of specific Location information items including
1022    Telephone, URL, Email, InstantMessaging, and RegionalCoverage.  Two
1023    approaches are possible. First, use the items directly under Organization to
1024    provide the current information and place Address, Country, and
1025    GeographicLocation in a Location element with a type of "current". Move all of
1026    the current information (from both Location and elements under Organization to a
1027    Location element with a type "superceded" when location information changes.
1028    This could also include a note providing its valid coverage dates. The
1029    organization entry would have a version update and the VersionRational would
1030    state that there were changes in the contact information.
1031
1032    Organization name changes can be handled in two ways. First if the organization
1033    is essentially the same simply with a new name and you wish to retain the ID. To
1034    retain the original ID, version the Organization entry, moving the old
1035    OrganizationName to Nickname and entering the new organization name along
1036    with a VersionRationale comment. If the organization has changed, for example
1037    merged or changed mission, you may wish to provide a new entry with a new ID.
1038    You can then use Relation to link it to the old ID and note the effective dates of
1039    the old and new organizations as well as descriptive and/or keyword
1040    classifications of the relationship.
1041
1042    Careful consideration of both the information you wish to retain and how it will be
1043    maintained will help you determine how you may want to control the use of
1044    options within OrganizationScheme. If this involves disallowing specific elements,
1045    making an element required, or providing it with a default value, the archive is
1046    advised to create a DDIProfile detailing this specialized usage. In addition, the
1047    archive may wish to develop controlled vocabularies to use in managing the
1048    information held in this scheme.
1049

1050    Entries for Individual are structured in much the same way as Organization.
1051    Individuals and Organizations can be nested as desired. Users may wish to
1052    decide to focus on either nesting or relation references as a primary structure for
1053    internal relationship definitions.

1054    ### *4.3   Life Cycle Information*

1055    Life cycle Information is a simple listing of important events in the life cycle of the
1056    study or group of studies. Each event uses an IdentifiableType as an extension
1057    base and  includes an EventType, a Date (this can be a simple date or a range),
1058    a reference to the organization or individual responsible for the event, and a
1059    description of the event. EventType is a CodeValueType which allows an
1060    organization to define specific codes relevant to their organization. DDI may
1061    publish a basic list of lifecycle events to provide a common list for DDI users. Life
1062    cycle events should be entered as they occur.
1063
```
1064    <r:LifecycleInformation>
1065    <r:LifecycleEvent isIdentfiable="true" id="LC_1">
1066            <r:EventType codeListID="DDI_Events_List" codeListAgency="DDI Alliance" >
1067            StudyDesign
1068            </r:EventType>
1069            <r:Date><r:SimpleDate>2007-04-01</r:SimpleDate></r:Date>
1070            <r:AgencyOrganizationReference>
1071            <r:Reference>
1072            <r:URN>
1073            urn:ddi:3_0:OrganizationScheme.Organization=MPC:MAINORG[1.0].MPC[1.0]
1074            </r:URN>
1075            </r:Reference>
1076            </r:AgencyOrganizationReference>
1077            <r:Description>
1078            Initial study proposal outlined and documentation  started.
1079            </r:Description>
1080    </r:LifecycleEvent>
1081    </r:LifecycleInformation>
```

1082    ## 5.0   Building and Documenting a Questionnaire

1083    Data collection contains the major components of the questionnaire's content
1084    and structure. In essence, a questionnaire contains questions and response
1085    options in an organized arrangement generally with instructions for the person
1086    answering either directly or through the interpretation of an interviewer. DDI 3.0
1087    distinguishes all of these parts primarily to support reuse through category
1088    schemes, coding schemes, and question banks. Separating these component
1089    parts allows for a generic description of instrument flow that can be captured
1090    from computer-aided survey systems or used to instruct such a system.
1091
1092    To enter information into this section of Data Collection, you need to think about
1093    a survey as a collection of its component parts. Some parts are persistent in
1094    relationship to the content of the question and others are related to the use of a

1095   question in an instrument. The following chart differentiates these major
1096   divisions. These are entered and assembled in life cycle order, selecting the
1097   questions needed to measure the desired concepts, determining the allowable
1098   responses, setting up question order, and adding related text to guide the
1099   respondent through the questionnaire and inform capture of the raw data.
1100

| Persistent Content of the Question | Use of the Question in an Instrument |
|---|---|
| Question text<br>• Simple text<br>• Dynamic text | Order and Routing<br>• Sequence<br>• Loops<br>• Skip patterns |
| Multiple part question [a question whose comparability is based on a particular question sub-series] | Universe |
| Response Domain<br>• Open<br>• Set categories<br>• Special types (date, time, etc.) | Analysis Unit |
| Definitional text [specific to the question as opposed to its use in a particular instrument] | Instructions<br>• Enumerator<br>• Respondent<br>• Coding (capturing raw data) |

1101
1102   This type of deconstruction is needed to allow for support of question banks that
1103   both reuse question text and associate multiple response domains with the same
1104   question.
1105
1106   Questions are structured to support capabilities for dynamic text, multi-language
1107   comparability, and alternative response domains. The response to the same
1108   question can often be captured in multiple ways, for example:
1109
1110        Question: How old are you?
1111        Response Domain 1:
1112          Numeric, maximum 3 characters
1113        Response Domain 2:
1114          Under 18 years
1115          18 to 64 years
1116          65 Years and older
1117        Response Domain 3:
1118          Under 5 years
1119          5 to 9 years
1120          10 to 14 years
1121          …continuing in 5 year cohorts
1122
1123   In addition, separating the flow logic, interviewer instructions, and identification of
1124   external materials used in the questioning process, allows DDI 3.0 to support a

1125 wide number of instrument types by providing generic instructions that can be
1126 interpreted by the individual instruments for presentation in multiple media
1127 formats.
1128
1129 Identifying these parts when looking at questionnaire can be more difficult.
1130
1131 IF LongIll=Yes THEN
1132      FOR i=1 to 6 DO
1133           IF (i=1) OR (More[i-1]=Yes) THEN
1134                *Records up to six long-standing illnesses*
1135                (IllsM[i])
1136                What (*else*) is the matter with you?
1137                INTERVIEWER: RECORD FULLY. PROBE FOR DETAIL. IF
1138                MORE THAN ONE MENTIONED, ENTER ONE HERE
1139                ONLY.
1140
1141                Text: Maximum 60 characters
1142
1143 This question starts with question routing information that is specific to the
1144 instrument. It would use the IfThenElse construct as its initial organizing
1145 structure. This is followed by a Loop construct (or possibly a repeat while) that
1146 repeats the question up to 6 times. The question itself is found in another
1147 IfThenElse which prompts the inclusion of the dynamic text (else) in the question
1148 text. The text "*Records up to six long-standing illnesses*" appears to be the
1149 purpose of the question, but is actually the reason for looping the question within
1150 this questionnaire. As such it should be included in control construct as an
1151 InteviewerInstructionRefereence. The text "INTERVIEWER: RECORD FULLY…"
1152 would also be an InterviewerInstruction although it may be considered part of the
1153 question or question construct rather than as part of the IfThenElse or Loop
1154 construct. The question text is dynamic;
1155
1156  <QuestionText>
1157       <LiteralText>What</LiteralText>
1158       <ConditionalText>
1159       <Expression>IF i>1 THEN: else</Expression>
1160       </CondidtionalText>
1161       <LiteralText>is the matter with you?</LiteralText>
1162  </QuestionText>
1163
1164 "Text: Maximum 60 characters" provides information on the type of
1165 ResponseDomain (text) and the maximumLength (60) of the response.
1166
1167 Another common occurrence in printed questionnaires is the inclusion of routing
1168 information next to response categories.
1169
1170 H10a TOILET FACILITIES: What type of toilet is used by the household?
1171

1172      1. W.C.
1173      2. Pit Latrine
1174      3. KVIP
1175      4. Bucket/Pan
1176      5. Toilet in another house (different house)      (Go to H11)
1177      6. Public Toilet (WC, KVIP, Pit, Pan etc.)      (Go to H11)
1178      7. No facilities (bush/beach/field)      (Go to H11)
1179      8. Other (specify) _____
1180

1181 In the above question there is both routing directions IF H10a > 4 AND H10a < 8
1182 THEN H11 ELSE H10b. Also category 8 requires a write in response. In coded
1183 responses this would involve the use of another IfThenElse requesting a text
1184 response to specify the Other category using a separate question.
1185

1186 Occasionally complex responses are used including a check box, fill-in number,
1187 and categories. For example:
1188

1189 48b how much is your regular monthly payment on all second or junior
1190 mortgages and all home equity loans on THIS property?
1191      Monthly amount—Dollars
1192      $_ _ , _ _ _ .00
1193

1194        OR
1195   _ No regular payment required
1196

1197 This is an example of a StructuredMixedResponse that contains:
1198

1199 <ResponseText><LiteralText>Monthly amount—Dollars</LiteralText></ResponseText>
1200 <NumericDomain>...</NumericDomain>
1201 <ResponseText><LiteralText>OR</LiteralText></ResponseText>
1202 <CategoryDomain>...</CategoryDomain>
1203

## 5.1   Question Construction

1204

1205 A basic question is an extention of VersionableType and contains a question text,
1206 question intent, response domain, concept references, and related visual aids.
1207 The presence of question intent allows you to use this structure for question
1208 development by capturing the intent of the question. This is also helpful when
1209 collecting information used as an indicator or when dynamic text changes or
1210 language differences result in varied wording of the question itself. The question
1211 text contains a Text element that serves as a container for a variety of text types
1212 and provides for a description and content. Note that questions do not contain
1213 display information for the content as this is specific to its use in an instrument.
1214

1215 Text types can currently be:
1216

| Text Type | Usage |
|---|---|
| Literal Text | The value is a static text string. |
| Conditional Text | Provides a condition on which the associated text varies (for example, gender). |

1217

1218 All questions have a designated response domain. This may be a reference to
1219 previously defined category schemes or coding schemes using the Response
1220 Domain Reference, or by defining the response domain within the question.
1221

| Response Domain Type | Usage |
|---|---|
| Text | Text content is used for open ended, non-numeric, or mixed response. They may be constrained by acceptable content range and/or by length. |
| Date Time | Structures a specified data and/or time content |
| Numeric | Provides a numeric type code, scale, decimal positions, start and end value of accepted range, and the allowed interval. |
| Code | References a defined CodeSheme. CodeSchemes are used by both questions and variables although the same scheme may not be used by a question and its resulting variable due to recoding. |
| Category | References a defined CategoryScheme. A category scheme is a list of valid responses that do not have code representations assigned to them. |
| Geographic | A structure for capturing required information from GPS systems and other means of geographic location definition used in a data collection instrument. |
| StructuredMixed | Allows for use of multiple response domain types plus DynamicText in order to structure complex response domains that contain multiple response types. |

1222

1223 CodeSchemes and CategorySchemes are defined in the logicalproduct and are
1224 used to describe the response domains of both questions and variables. In
1225 general, a category scheme is organized and assigned representational codes by
1226 a code scheme prior to its use in a variable such as code 0 equals Male and
1227 code 1 equals Female. A question may or may not use a coded category to
1228 capture a response. Coding is assigned when transferring a response to the raw
1229 data file.
1230

1231 Category schemes are used when no code is provided in the instrument for the
1232 selected answer and coding instructions provide information on how the selected
1233 response is captured in the raw data. This is most common in paper-based
1234 questionnaires, for example:
1235

1236      Question:
1237       What is your marital status?
1238      Response Domain:

1239             ◯ Married

1240             ◯ Single, never married

1241             ◯ Widowed

1242             ◯ Divorced

1243
1244  Question Schemes define a set of questions used in a study. They may include
1245  inline descriptions and/or references to external questions from question banks.
1246
1247  MultipleQuestions are those that require explicit responses to a set of
1248  subquestions. For example:
1249
1250  Q10. Read through the following list of candidates and for each candidate
1251  indicate how familiar you are with his position on free trade.

1252             1    2    3    4
1253      Bill Clinton  O    O    O    O
1254      Bob Dole    O    O    O    O
1255      Ross Perot  O    O    O    O
1256
1257      Where        1 = Very familiar
1258                  2 = Somewhat familiar
1259                  3 = Not familiar
1260                  4 = Do not know of this candidate
1261
1262  The sub-questions consist of each candidate's name as the question text and the
1263  CodeDomain.
1264
1265  Sequencing of response categories, subquestions of a multiple question, and
1266  questions within a series can be defined using an element provided in
1267  QuestionConstruct (ResponseSequence), MultipleQuestion
1268  (SubQuestionSequence), and Seqence (ConstructSequence). The user can
1269  designate the default of InOrderOfAppearance, Random, Rotate (rotates through
1270  a sequence), Other (defined in AlternateSequence).

## 5.2   Control Constructs and Instrument

1272  The purpose of the DDI instrument element is to record the flow of a
1273  questionnaire, its use of questions, and additional component parts. Instrument is
1274  a maintainable object and has a MaintainableType. The documenter can define
1275  instrument types and provide a type code for the instrument as a whole. Software
1276  used to collect data is identified here along with a reference to other objects
1277  associated with the instrument such as a physical or image copy of the

1278  questionnaire. Other than this general instrument classification information,
1279  Instrument is composed of a series of Control Constructs nested inside a single
1280  master Control Construct. Generally the top level Control Construct would be
1281  Sequence, implying a start and end of the instrument as the Sequence tag is
1282  opened and closed.
1283

| Control Construct | Use |
| --- | --- |
| IfThenElse | Provides a condition to trigger the Then clause and alternative Else action. Then and Else are both Nested Constructs. |
| RepeatUntil | Provides an Until Condition and a construct that continues until the condition is met. |
| RepeatWhile | Provides an While Condition and a construct that continues while the condition is true. |
| Loop | Identifies the Loop Variable, sets its initial value, the condition that must be met to end the loop, and step (increment) value for the loop variable and the control construct that is to continue until the end of the loop is met. |
| Sequence | Defines a sequence of control constructs reflecting the flow of the instrument or a subsection of an instrument. |
| ComputationItem | Used to assign a code to a variable. |
| StatementItem | Statement type provides for the inclusion of additional text such as pretext or posttext. |
| QuestionConstruct | Provides for the insertion of a question into the instrument. In addition to the question content, the sequence includes ability to designate a response unit, analysis unit, and to reference a universe. |

1284
1285  A ControlConstructScheme lists all the constructs found within the instrument as
1286  individual items. Constructs that define the organization of other constructs such
1287  as Sequence or IfThenElse, Loop, etc. do not include their subordinate
1288  constructs inline but by reference. The easiest approach is to make a list of all
1289  Questionconstructs, StatementItems, and ComputationItems so that they are
1290  ready to include by reference. Then review your flow logic, identifying routing
1291  patterns, sequences of QuestionConstructs and/or StatementItems and begin
1292  defining these from the lowest level of the nesting out until you have a construct
1293  (generally a sequence) that provides the start and end point for the full
1294  instrument. This is the ControlConstruct that the Instrument will point to.
1295
1296  For example if you had a questionnaire with a single skip pattern, 6 questions,
1297  and 2 statements you might have the following contents in your
1298  ControlConstructScheme:
1299
1300  Question 1

| | | |
|---|---|---|
| 1301 | Question 2 | if 4 GO TO Statement 2, Question 4 |
| 1302 | | Else GO TO Statement 1, Question 3, Question 4 |
| 1303 | Question 5 | |
| 1304 | Question 6 | |
| 1305 | | |

```
1306    <ControlConstructScheme isMaintainable="true" id="CC1">
1307          <QuestionConstruct isVersionable="true" id="QC1">...</>
1308          <QuestionConstruct isVersionable="true" id="QC2">...</>
1309          <QuestionConstruct isVersionable="true" id="QC3">...</>
1310          <QuestionConstruct isVersionable="true" id="QC4">...</>
1311          <QuestionConstruct isVersionable="true" id="QC5">...</>
1312          <QuestionConstruct isVersionable="true" id="QC6">...</>
1313          <StatementItem isVersionable="true" id="SI1">...</>
1314          <StatementItem isVersionable="true" id="SI2">...</>
1315          <Sequence isVersionable="true" id="SQ3">
1316                <ControlConstructReference isReference="true">
1317                      <ID>SI2</ID>
1318                </ControlConstructReference>
1319                <ControlConstructReference isReference="true">
1320                      <ID>QC4</ID>
1321                </ControlConstructReference>
1322          </Sequence>
1323          <IfThenElse isVersionable ="true" id="IF1">
1324                <IfCondition><Code>QC2 = 4</Code></IfCondition>
1325                <ThenConstructReference isReference="true">
1326                      <ID>SQ3</ID>
1327                </ThenConstructReference>
1328                <ElseConstructReference isReference="true">
1329                      <ID>SQ2</ID>
1330                </ElseConstructReference>
1331          </IfThenElse>
1332          <Sequence isVersionable="true" id="SQ2">
1333                <ControlConstructReference isReference="true">
1334                      <ID>SI1</ID>
1335                </ControlConstructReference>
1336                <ControlConstructReference isReference="true">
1337                      <ID>QC3</ID>
1338                </ControlConstructReference>
1339                <ControlConstructReference isReference="true">
1340                      <ID>QC4</ID>
1341                </ControlConstructReference>
1342          </Sequence>
1343          <Sequence isVersionable="true" id="SQ1">
1344                <ControlConstructReference isReference="true">
1345                      <ID>QC1</ID>
1346                </ControlConstructReference>
```

```
1347            <ControlConstructReference isReference="true">
1348                    <ID>QC2</ID>
1349            </ControlConstructReference>
1350            <ControlConstructReference isReference="true">
1351                    <ID>IF1</ID>
1352            </ControlConstructReference>
1353            <ControlConstructReference isReference="true">
1354                    <ID>QC5</ID>
1355            </ControlConstructReference>
1356            <ControlConstructReference isReference="true">
1357                    <ID>QC6</ID>
1358            </ControlConstructReference>
1359        </Sequence>
1360 </ControlConstructScheme>
1361
```

## 6.0  Data Processing

1363 Data processing takes place at various points in the life cycle. The first such point
1364 concerns instructions for translating the response to the question into the raw
1365 data file. This is normally a direct capture. However, in the case of paper
1366 questionnaires this may include attaching codes to response categories or
1367 incorporating information not collected from the questionnaire. This could be
1368 geographic information, identifiers, recoded, or derived items. In addition, general
1369 instructions may address the handling of non-response, illegal multiple
1370 responses, or other common coding and processing issues.
1371
1372 The specific areas of information captured in DataProcessing include control
1373 operations, cleaning operations, weighting factors, data appraisal, and coding.
1374 Both ControlOperation and CleaningOperation contain a repeatable description
1375 field and an AgencyOrganizationReference to an organization or individual
1376 described in Organization. When entering this information separate out different
1377 operations or different steps. ControlOperation and CleaningOperations as well
1378 as the Description elements within them are unbounded. While these currently
1379 are not identifiable elements, keeping discrete activities separate is in line with
1380 the general philosophy of DDI 3.0.
1381
1382 Weighting factors are identified structured strings allowing for specific weight
1383 factors to be listed discretely and referenced. Although this allows for structured
1384 statements, consider separating descriptions of weighting factors (which may be
1385 structured) from any standard weights that are used as factors in category or
1386 variable statistics. You will want to refer to this as a number that was used in
1387 calculation.
1388
1389 Data Appraisal Information provides response rate (string), sampling error
1390 (structured string), and OtherAppraisalProcess (structured string). Once again, all

1391  of these are repeatable, and maintaining clear and discrete information on
1392  processes is important.

## 6.1   Coding

1394   There are two different types of Coding describing general processes that are
1395  applied to a wide range of variables and specific generation instructions. All
1396  Coding elements are an extension of IdentifiableType so they can be referenced
1397  by ControlConstructs and Variables. Coding contains either GeneralInstruction or
1398  GenerationInstruction.
1399
1400  A general instruction pertains to coding operations such as the uniform handling
1401  of non-response or general imputation instructions. It contains a required
1402  repeatable Description and Command structure. The first is a structured string
1403  and intended to be human-readable. Command provides for a human-readable
1404  CommandText, an optional repeatable reference to a CommandFile, and an
1405  optional StructuredCommand that allows for inserting extentions to provide
1406  structured language for external namespaces such as MathML. When creating
1407  GeneralInstructions make sure that the relationship between the Description and
1408  Command is clear. Do not mix several Descriptions and/or Commands pertaining
1409  to multiple instructions in a single GeneralInstruction. Create a separate
1410  GeneralInstruction for each. GeneralInstruction also contains an element and
1411  attribute pair IsOverride and @isOverride. When @isOverride is changed from its
1412  default value of "false" to "true", the element IsOverride must be completed. This
1413  element provides the ID of the GeneralInstruction that is being overridden.
1414
1415  GenerationInstruction is more complex as it deals with specific recodes and
1416  derivation instructions. An attribute, isDerived, indicates whether this is a simple
1417  recode or a more complex derivation instruction. GenerationInstruction In
1418  addition to the Description and Command structures found in GeneralInstruction
1419  it provides identification of Questions, Variables, and ControlConstructs used in
1420  the generation command. You are able to assign a mnemonic to the Question or
1421  Variable for use in the command equation.  It also contains a special Aggregation
1422  description which provides the aggregation method and identifies the
1423  Independent and Dependent variables required for the computation. This
1424  structure is also available in CodeMap for use in defining recodes from the
1425  Source to the Target structure.

# 7.0   Creating a Basic Data Dictionary

1427  A basic data dictionary in its simplest form consists of a description of variable
1428  contents, information on the universe, the concept represented by the variable,
1429  the measurement unit, analysis unit, the number of respondents, and the location
1430  of the data field in the physical data store. DDI 3.0 separates these parts into
1431  separate category schemes, code schemes, variable descriptions, and physical
1432  location.

## 7.1   *Category Schemes*

1433

1434 Category schemes should be started early in the process. Category schemes are
1435 used by questions as response domains and by code schemes in preparation for
1436 use in variables. A category scheme can include a set of related categories or
1437 include all the categories used in the study. Note that category schemes are
1438 maintainable and you should consider creating separate category schemes for
1439 subsets of categories that may be reused in other studies. This is especially true
1440 if your questionnaire uses un-coded categories. One group of categories you
1441 may wish to construct is one including non-response categories used in the
1442 study. Any category used in the study should only be listed once. Reuse of a
1443 category in multiple code schemes implies comparability.
1444
1445 As a maintainable object, Category Schemes are an extention of a
1446 MaintainableType and contain descriptive information about the scheme (label
1447 and description) to provide information on what the scheme contains. Its content
1448 is made up of a set of category descriptions which may or may not be organized
1449 into category groups. A category is an extension of VersionableType to allow for
1450 updates and changes over time, such as changes in a label or description. The
1451 label is repeatable to handle multiple languages. This is the actual category term
1452 as used in the question response domain or variable representation. The
1453 description is repeatable to handle multiple languages. When determining the
1454 comparability of categories, base the comparison on the description not on the
1455 label. For example, the category "Chemist" in British occupation codes has the
1456 same description and is comparable to the American occupation "Pharmacist"..
1457 In addition, the object Generation allows for description of the command used to
1458 generate the contents of this category, for example "Other Income" equaling total
1459 income minus wage/salary income.

## 7.2   *Code Schemes*

1460

1461 Code schemes organize categories from one or more CategorySchemes and
1462 provide the code representation for the category as it is found in the question or
1463 variable. A code scheme may be flat or hierarchical and the hierarchy can be
1464 regular (each branch having the same number of levels) or irregular (the number
1465 of levels varies by branch). A code within a code scheme is identified by its
1466 unique code value. A level or specific code may or may not have data associated
1467 with it. A code without associated data acts as a category group as found in
1468 earlier DDI structures. Its sole purpose is to group a subset of categories/codes
1469 and provide a grouping label through reference to a category.. A CodeScheme
1470 contains only a single code scheme and should provide all legitimate categories.
1471 If a standard codeScheme is used for non-response, a single CodeScheme can
1472 describe that and be included in each relevant codeScheme by reference.
1473
1474 To build a CodeScheme you must have completed the relevant
1475 CategorySchemes. A CodeScheme has a MaintainableType, Label, and
1476 Description. It may reference one or more existing CodeSchemes. For example,
1477 if there is an existing CodeScheme for Gender (0 = Male, 1 = Female) and

1478 another for Non-response (9 = No response), a single CodeScheme could be
1479 created to include these three valid responses, 0 = Male, 1 = Female, and 9 = No
1480 response. If all or a majority of the categories come from a single
1481 CategoryScheme, it can be declared once, thereby creating a default value for
1482 the maintainable object of category references allowing subsequent category
1483 reference to list only the category's ID and Version information. This can be
1484 overridden by providing a complete URN at the category reference level. If the
1485 CodeScheme is hierarchical, you must indicate if it is regular or irregular. Flat
1486 CodeSchemes are the default and do not require this element. Hierarchical
1487 CodeSchemes must have their levels described so that they can be referenced
1488 by Variables using a limited number of levels and to understand the relationship
1489 between the levels and their contents.
1490
1491 A level has a name, description, relationship type, and interval. The relationship
1492 type describes the relationship of the categories contained by the level. They can
1493 be nominal (no implied order), ordinal (ordered as presented in the description),
1494 interval (both ordered and with a consistent interval between categories), ratio
1495 (ordered and with a consistent interval ratio), or continuous (either interval or
1496 ratio, use when unsure of the interval type) . If the categories have an interval
1497 relationship, provide the anchor (base value for the first category) and increment
1498 value of the interval. This information is provided for each level as this may vary
1499 by level, for instance, a single hierarchy may have ordinal or interval relationships
1500 at upper levels and nominal at the lowest level.
1501
1502 Codes are then listed containing a reference to the category being represented
1503 and a code value, plus a nested code to allow for building hierarchies. A code
1504 has two attributes, a level number (optional) to indicate the level of the coded
1505 category, and isDiscrete. The field isDiscrete has a default value of "true". Set
1506 this attribute to "false" if it has a subordinate level.
1507

1508 ## 7.3   Describing Variables

1509 A variable is part of a variable scheme and is made up of the following
1510 components:
1511

| Element | Usage |
|---|---|
| VersionableType extension base | This is how the variable is referenced by other objects in the instance |
| Name | Contains an optional Name for the variable such as a mnemonic. This is available in all identifiable elements but is most commonly used in Variable. Note that it is repeatable for language and geographic differences. |
| Label | This is a short human-readable label description of the variable. |
| Variable Definition | Additional textual description of the variable |

| | which may be used to provide extended detail concerning the variable. |
|---|---|
| Universe Reference | Universe is described in ConceptualComponent. This allows the universe of the variable to be seen in context of the full universe structure and provides comparability for two variables using the same universe. A universe is assumed to be the fullest universe of the study if it is not provided but given the complexity of multiple data products being produced from a single study, it is recommended that a universe be explicitly identified. Variables used as component parts (attributes or dimensions) of NCubes do not require a universe as this is defined in the NCube. |
| Concept Reference | Similar to Universe in terms of concepts being described within the ConceptualComponent. All variables must have one and only one concept declared. |
| Question Reference | References to all questions as expressed in an instrument used to determine the value of this variable. |
| Embargo Reference | References access restrictions for this variable. |
| Response Unit | Who provided the information for this variable. This may be the same as the response unit of the question as used in the instrument. |
| Analysis Unit | Describes who or what this variable is an attribute of (who or what is described). |
| isTemporal isGeographic isWeight | Three Boolean attributes all with the default value of "false". If true, reset the value of appropriate attribute to "true". |
| Representation: With optional attributes of measurementUnit, aggregationMethod, additivity | Describes how the variable is represented in the data file according to one of the following. Attributes are optional as this information is listed in the NCube when the variable is used as a dimension or attribute of an NCube. |
| WeightVariableReference | References the variable to be used with this variable. |
| StandardWeightReference | References the weight factor described in Data Processing. |
| ConcatenatedValue | Describes 2 or more variables, that when concatenated provide the value of this variable. Used primarily to create a virtual variable such as a multipart link whose parts are described by separate variables. |

| CodingInstructionReference | References the specific coding instruction in Coding used to create this variable. This may be a recode or derivation instruction. |
| ValueRepresentation | Provides the actual representation used by the variable. This is an abstract object that acts as head of a substitution group. |

1512
1513
1514 Some of these elements may have been compiled at earlier stages of the study.
1515 For example, the concept and universe structures, questions, a number of
1516 category and possibly code schemes, and some of the recoding or derivation
1517 process will already exist. The first step is to complete those component parts.
1518
1519 The Value Representations available are described below. Note that
1520 CategorySchemes and CodeSchemes should be created first. CodeSchemes
1521 should describe the full structure of the code. A variable can use all or parts of
1522 the CodeScheme thereby retaining the relationships among the various parts or
1523 levels of the structure. Each of these types extend the contents of
1524 VariableRepresentation noted in the above table.
1525

## 7.3.1  Text

1526
1527 TextRepresentationType allows you to specify the minimum and maximum length
1528 of the text content and to provide a regular expression that can be used to restrict
1529 the content of the text string. For example, a United States ZIP Code is a text
1530 string because the leading zeros carry significance. As a text string a five digit
1531 ZIP Code would have a minLength of 5, a maxLength of 5, and a regular
1532 expression of [0-9]*. Please note that text letters that represent categories labels
1533 are described with CodeScheme. The value of a code may contain any
1534 character.
1535

## 7.3.2  Date/Time

1536
1537 DateTimeRepresentation requires the use of one of the following types
1538 corresponding to the W3C XML schema xs: datatype.
1539

| Code | W3C XML schema xs datatype |
|---|---|
| DateTime | [xs:dateTime] Contains both the date and time as <date>T<time> |
| Date | [xs:date] Contains the full date from the Gregorian calender YYYY-MM-DD unless an alternative format is provided |
| Time | [xs:time] Contains the full time on a 24-hour clock system unless alternative format is provided. hh:mm:ss. Precision can be dropped resulting in hh:mm or hh.  A time zone can be added <time>Z using the standard time zone designation +-hh:mm or +-hh |
| Year | [xs:gYear] Contains the 4 digit year YYYY |

| Month | [xs:gMonth] Contains the 2 digit month MM |
| --- | --- |
| Day | [xs:gDay] Contains the 2 digit day DD |
| MonthDay | [xs:gMonthDay] Contains the 2 digit month followed by the 2 digit day as MM-DD unless an alternative format is provided |
| YearMonth | [xs:gYearMonth] Contains the 4 digit year followed by the 2 digit month as YYYY-MM unless an alternative format is provided |
| Duration | [xs:duration] Provides a duration of time represented by one of the following formats (specific format must be declared) PnnYnnMnnDTnnHnnMnnS where n is replaced with the number of unit types for example "P3Y6M4DT12H30M0S" defines "a period of three years, six months, four days, twelve hours, thirty minutes, and zero seconds". Elements may be omitted if their value is zero. T is used to separate date and time elements so that P3M is 3 months and PT3M is three minutes. Alternative format P<date>T<time> "P0003-06-04T12:30:00". |
| Timespan | This is not allowed as a date type when describing an NCube dimension as it represents two dimensions. Complex structure containing <start>/<end>, <start>/<duration>, or <duration>/<end>. Start and end can follow any of the designated datetime structures and should be declared in format. <start>/<end> example: "2007-03-01T13:00:00/2008-05-11T15:30:00" <start>/<duration> example:  "2007-03-01T13:00:00/P1Y2M10DT2H30M" <duration>/<end> example "P1Y2M10DT2H30M/2008-05-11T15:30:00" For <start>/<end> expressions, if any elment are missing from the end valude, they are assumed to be the same as for the start value including the time zone if used. For example a 2 hour meeting "2007-12-14T13:30/15:30". |

1540

1541 An optional format element can define an alternative format for the data field
1542 such as MM-DD-YY. The default is the W3C format.

1543

1544 **7.3.3  Numeric**
1545 Note that numeric should not be used for numbers that represent categories (1 =
1546 Male, 2 = Female). The use of a numeric representation suggests a form of count
1547 (6 years, 20000 Euros, 2 children). Numeric representation provides a start and
1548 end value for the range, scale, the number of decimal positions, interval, and a
1549 numeric type code. The type code is required and is represented by one of the
1550 following.

1551

| Code | W3C XML schema xs datatype |
| --- | --- |
| BigInteger | [xs:integer ]An integer of unlimited size. An integer datatype |

| | corresponding to W3C XML Schema's xs:integer datatype. |
|---|---|
| Integer | ]xs:int] An integer number can hold a whole number, but no fraction. Integers may be either signed (allowing negative values) or unsigned (nonnegative values only). An integer datatype corresponding to W3C XML Schema's xs:int datatype. |
| Long | ]xs:long] An integer of up to 32 bits in size (corresponding to an unsigned range of 0 to 4,294,967,295 or a signed range of -2,147,483,648 to +2,147,483,647). A numeric datatype corresponding to W3C XML Schema's xs:long datatype. |
| Short | ]xs:short] An integer of up to 16 bits in size (corresponding to an unsigned range of 0 to 65,535 or a signed range of -32,768 to +32,767), A numeric datatype corresponding to W3C XML Schema's xs:short datatype. |
| Decimal | ]xs:decimal] A real number (allows fractions expressed as decimals). A numeric datatype corresponding to W3C XML Schema's xs:decimal datatype. |
| Float | ]xs:float] Real numbers that may be stored in scientific notation (example: 20.0005, 99.9, -5000.12, 6.02e23). A numeric datatype corresponding to W3C XML Schema's xs:float datatype. |
| Double | ]xs:double] Float of up to 32 bits. A numeric datatype corresponding to W3C XML Schema's xs:double datatype. |
| Count | Ordinal number of objects in a finite set, discrete. A simple incrementing Integer type. The isSequence facet must be set to true, and the interval facet must be set to "1". |
| Incremental | A value that is continuous and infinite can be interval or ratio. This value indicates that the value increments according to the value provided in the interval facet, and has a true value for the isSequence facet. |

1552
1553
1554 **7.3.4  Code**
1555 The value representation for code allows a single CodeScheme to be applied in a
1556 number of ways. You can define which portions of the CodeScheme are being
1557 used in a particular variable. The default is to include all codes described in the
1558 codeScheme. In the case of a hierarchical codeScheme, this means that the
1559 response can include all levels of the hierarchy. Alternatively, the following
1560 objects can be used to constrain what is included in the variable.
1561

| IncludeLevel | Identify specific levels to be included in the variable |
|---|---|
| IncludedCode | Reference to included codes |
| DataExistence | Include the most discrete items only (this allows inclusion of all codes designated as the most discrete which could be a level or in the case of an irregular hierarchy, the rightmost |

| codes in each branch of the tree's hierarchy |
| --- |

1562
1563   Example:
1564
1565   Irregular Hierarchy
1566   1  Metals (level 1)
1567   2        Iron  (level 2)
1568   3              Bar Iron (level 3, isDiscrete="true")
1569   4              Caste Iron (level 3, isDiscrete="true")
1570   5        Copper (level 2, isDiscrete="true")
1571   6  Non-Metals (level 1, isDiscrete="true")
1572
1573   Variable 1 – using IncludeLevel
1574     Includes Level 1 and Level 2
1575     Valid responses are 1, 2, 5, 6
1576
1577   Variable 2 – using IncludedCode
1578     Includes valid responses 3, 4, 5
1579
1580   Variable 3 – using DataExistence
1581     Includes valid responses 3, 4, 5, 6
1582
1583   Variable 4 – no constraints noted
1584     Includes valid responses 1-6

## 1585   *7.4    Data Relationships*

1586   The Data Relationship portion of a logical product describes the logical records
1587   described in terms of their coverage, unique identifiers, and interrelationships.
1588   This section is meant to reflect the information often found in the "How to use this
1589   file" type sections of a traditional codebook. All data sets have one or more
1590   logical record types. Note that this is referring to the logical structure of the
1591   records not their physical layout which may be hierarchical, rectangular, or held
1592   in a relational data set. This section is focused on what variables are available to
1593   link data together and to assist the user in identifying unique cases with a record
1594   type. At a minimum, the human-readable description should be completed.
1595   Ideally the detailed information should be completed to support machine-
1596   actionable exploration of the data set.
1597
1598   The Data Relationship section contains three main parts. First is a human-
1599   readable description of the logical record types contained in the data set and the
1600   relationships among these record types. The second part describes each logical
1601   record type. The record type descriptions are Identifiable so that they can be
1602   referenced by the physical data structure.
1603
1604   The variable value reference provides the identification of the variable that
1605   differentiates one record type from another and the value of this variable for the

1606    record type being described. In files with a single record type this would not be
1607    used. However, most files with more than one record type have a variable such
1608    as RecordType that provides a value such as "H" for household record and "P"
1609    for person record. Older data files may have an identifier for the first record with
1610    subsequent records (records of another type nested within them) simply being
1611    those records that do not have this value. The attribute hasLocator must indicate
1612    whether or not the record has such a variable.
1613
1614    Support for multiple parts allows for the identification of any variable that provides
1615    information that can be use to identify which segment of a logical record you are
1616    dealing with. This is a common practice for long logical records that are being
1617    stored within the constraints of a particular storage structure (Excel files have a
1618    record length, old SPSS and SAS packages had a maximum record length that
1619    could be supported). Support for multiple records is often found in elements such
1620    as LogicalRecordPartNumber, which appears in the 1990 and 2000 U.S. Census
1621    Standard Summary Files. Other older files have record segments without support
1622    for multiple parts, and identification is based solely on record order within the file.
1623    This section provides both the variable used to identify a record part and the
1624    values available.
1625
1626    Case specification provides for one or more means of identifying a unique case
1627    within a record type. For example, a record within its original file structure may
1628    have a unique record ID variable, often used to link it to other records. However,
1629    there may also be other identifiers, either individual fields or combinations of
1630    fields, that allow for case identification. A common example of this is the
1631    geographic codes of aggregate data files. Case specification allows for
1632    alternative identification for various case types. For example, in a geographic
1633    case file, a County level record may require one set of variables (state and
1634    county codes) to identify it while a Place requires another set (state and place).
1635    Case specification allows for different sets of identifier variables to be declared
1636    for different values of a parent identifier.
1637
1638    The last part of the record description is a listing of the variables belonging to the
1639    logical record. This is a special type of Variable Group and should be declared
1640    within the Data Relationship rather than as a standard Variable Group.
1641
1642    The third section of Data Relationship deals with record relationships. These are
1643    declared in a pair-wise fashion to account for all supported relationships. This is
1644    done by defining the variable and value of the linking variable in the source
1645    record and its matching point in the target record. The relationship is defined
1646    through the VariableValueRelationshipType attribute as "parent", "child", or
1647    "sibling". The value relationship by default is "equal", but can be set to
1648    "GreaterThan", "LessThan", "GreaterThanOrEqual", LessthanOrEqual", or
1649    "NotEqual". Multiple variable links should be handled by creating a variable with a
1650    concatenated value and using it as the link variable. This relationship is

1651  identifiable so that it can be referenced by the physical data product, reducing
1652  repetition of the information.
1653
1654  Note that if you are creating links between multiple logical product sections, you
1655  can put this information in a single logical product (defining all the relationships in
1656  the collection of logical products), or create a logical product that contains just
1657  the data relationship between all other logical products. This is useful when
1658  describing links between files in a longitudinal series, where there is a link
1659  between a person record occurring in one year and that person's record in
1660  another year or wave.
1661

## 7.5 NCubes

1662

1663  NCubes are used to describe the data matrices created through crosstabulation
1664  and aggregation of microdata. An NCube represents a matrix where each cell
1665  intersects each dimension of the matrix in one and only one location. For
1666  example, AGE by SEX by COUNTRY OF RESIDENCE. A cell has an "address"
1667  which provides its intersect code for each dimension. In this case, "5,1,6" would
1668  be the fifth value of AGE, the first value of SEX, and the sixth value of COUNTRY
1669  OF RESIDENCE.
1670
1671  An NCube is an extention of  a VersionableType, and contains Label,
1672  Description, a Universe Reference, ImputationReference, ResponseUnit, and
1673  AnalysisUnit. In this, it is very similar to a variable. An NCube must be
1674  constructed of one or more dimensions which are described as variables. The
1675  dimensions identify both a variable and a "rank" order so that the coordinate
1676  address will be clear. In the above example AGE has a rank of "1", SEX a rank of
1677  "2", and COUNTRY OF RESIDENCE a rank of "3". Additional attributes can be
1678  attached to all or part of the NCube. These could be cell level suppression flags
1679  described by a variable, footnotes or source notes, or whatever is required. By
1680  describing regions of the NCube with Coordinate Groups, one or more attributes
1681  can be attached to a cell, a dimension or dimension value, or a specific sub
1682  region of the NCube defined by the intersect values. A common use of the
1683  attribute within the NCube description is to identify cells that contain no data by
1684  definition. Many NCubes are created through the process of cross-tabulation and
1685  then collapsed when published to save storage or printing space. For example,
1686  the U.S. Census has a table of Number of persons per household by Household
1687  type where Household type contains Non-Family and Family Households. The
1688  number of person's categories run from 1 to 9 or more. This table is usually
1689  published as a one dimensional matrix Number of persons per household by
1690  household type.
1691

| LABELS | Cell coordinates |
|---|---|
| Nonfamily Household: | [no data used a category group label] |
| 1 | 1 |
| 2 | 2 |

| | |
|---|---|
| 3 | 3 |
| 4 | 4 |
| 5 | 5 |
| 6 | 6 |
| 7 | 7 |
| 8 | 8 |
| 9 or more | 9 |
| Family Household: | [no data used a category group label] |
| 2 | 10 |
| 3 | 11 |
| 4 | 12 |
| 5 | 13 |
| 6 | 14 |
| 7 | 15 |
| 8 | 16 |
| 9 or more | 17 |

1692
1693  Note that 1 person Family Household is missing. By definition a "Family" is two or
1694  more persons so this cell would always be zero or have no data listed.
1695
1696  The layout of the NCube structure in DDI 3.0 allows for the following table
1697  structure and attributes description.
1698

| LABELS | Dimension 1: | |
|---|---|---|
| Dimension 2: | Nonfamily Household | Family Household |
| 1 | 1,1 | 2,1 |
| 2 | 1,2 | 2,2 |
| 3 | 1,3 | 2,3 |
| 4 | 1,4 | 2,4 |
| 5 | 1,5 | 2,5 |
| 6 | 1,6 | 2,6 |
| 7 | 1,7 | 2,7 |
| 8 | 1,8 | 2,8 |
| 9 or more | 1,9 | 2,9 |

1699
1700  The attribute would reference a Variable that declares it's sole value as blank
1701  with the label "Definition as a Family Household requires a minimum of 2 people
1702  in the household". The the value is attached at the CoordinateGroup which
1703  defines the CoordinateGroup as any cell with the dimension 1 value of 2 and
1704  dimension 2 value of 1 (Family Household, 1 Person).
1705
1706  Attribute values can be declared within the metadata if they are consistent for all
1707  instances of the NCube (such as footnotes or definitional values) or can be
1708  stored in the data file and attached to the cell in the RecordLayout definition.
1709

1710 The Measure of an NCube must be described by a variable in order to capture its
1711 universe and various aspects of the measure. For example, a measure of an
1712 NCube Persons Per Housing Unit may have a dimension denoting the number of
1713 persons per housing unit and the measure is a count of housing units falling into
1714 each definition. Alternatively, it may be the percentage of housing units with the
1715 designated number of people. In order to capture both the dimension and the
1716 measure, you need to have two unique definitions and their associated concepts.
1717 Some measures such as percentages require definition of the values that are
1718 used for the independent and dependent variables used to calculate the cell
1719 content. For example a percent could be the cell's percent of the universe or the
1720 cell percent of one of its dimensions. This can be handled one of two ways. First
1721 a separate variable can be created for each percentage or other aggregation
1722 using the specific variables to describe its calculation in GenerationInstruction. A
1723 variable pointing to this GenerationInstruction would be used for the measure.
1724 Alternatively a generic variable "Percentage of Housing Units" can be created
1725 with a GenerationInstruction that states that the value is calculated by dividing
1726 the dependent variable divided by the independent variable. The
1727 AggregationDefinition within Measure then declares which dimension(s) are used
1728 as the Independent variable and which are used for the Dependent variable. It
1729 can also note that the universe of the NCube is the independent variable. For
1730 example, when the NCube was Age by Sex and the Percent was provided for the
1731 each cell count divided by the total population of the NCube. If the percentage
1732 was the percent of 2 year olds who were female, Age would be the Dependent
1733 Variable and Sex would be the Independent variable.
1734
1735 Note that an NCube and a table are not synonymous.  A table may consist of a
1736 number of NCubes that share a common dimension, such as AGE by COUNTY
1737 OF RESIDENCE followed by SEX by COUNTRY OF RESIDENCE, where there
1738 is no intercept between AGE and SEX (for example no "Males 5 to 9 years of
1739 age").
1740

| Table 1 | Sex | | Age | | |
|---------|------|--------|----------|----------|------------|
|         | Male | Female | Under 18 | 18 to 64 | 65 or more |
| Region 1 | | | | | |
| Region 2 | | | | | |
| Region 3 | | | | | |
| Region 4 | | | | | |

1741
1742 In Table 1 there are two NCubes Sex by Region and Age by Region. By
1743 definition the cells of an NCube must intersect each dimension of the NCube at
1744 one and only one point. None of the cells in the above table intersect both the
1745 Age dimension and the Sex dimension. They share a common dimension of
1746 Region and will each point to the same Variable to describe this dimension.

## 1747 8.0 Physical Data Product

1748 Physical Data Product was designed to describe the physical storage structure of
1749 data including the breakout of logical data records into multiple physical parts, a
1750 link to the logical record being stored, and the specific structure of each record.
1751 The general structure of this section is split into describing the Gross structure of
1752 the physical records (what logical record is contained and how physical records
1753 should be linked), and the details of the RecordLayout (this may include data
1754 inline for some structures). Note that any logical product may have more than
1755 one physical storage structure. Frequently, files are stored in multiple formats, for
1756 example, a fixed format archive copy and a database copy that supports a
1757 specific use of the data file. In turn, each physical record description can
1758 represent multiple physical files of data. These files of data may represent
1759 different record types (different RecordLayout structures) within the dataset or
1760 different subsets of records within the full dataset (multiple PhysicalInstances)
1761 such as a datafile containing  just the records for the state of Arizona or just the
1762 records for females in the study.
1763
1764 Data can be stored in a wide and growing range of physical data structures. The
1765 common archival structures of fixed format and delimited files are described very
1766 much as they were in earlier versions of DDI. DDI 3.0 has tried to facilitate the
1767 management of the same data held in multiple storage structures and multiple
1768 datafiles. A single logical record of data may be storage as a single record string,
1769 broken into physical segments to accommodate line length limitations in some
1770 software, stored with all segments in the same file or different files (like relational
1771 data files), use a variety of proprietary and nonproprietary data structures, and be
1772 divided in multiple files of record subsets to facilitate processing. Within this,
1773 numerous pieces of information regarding each dataitem need to be recorded but
1774 much of it is repetitious and could be provided as default values. In order to
1775 handle this complexity, DDI 3.0 approaches the description of physical storage
1776 as follows:
1777

| Major Component | Use |
| --- | --- |
| Common Schema Components:<br>• MaintainableType<br>• OtherMaterial<br>• Note | Users may choose to maintain a single PhysicalDataProduct which describes their holding. Whether representing a single StudyUnit or a whole collection, OtherMaterial related to physical processing, analysis tools, systems information or other related materials or notes can be stored here. |
| PhysicalStructureScheme:<br>• MaintainableType<br>• Description<br>• Label | A separately maintainable scheme, PhysicalStructureScheme is basically a listing of individual PhysicalStructure descriptions included inline or through inclusion |

| | |
|---|---|
| • PhysicalStructureReference<br>• PhysicalStructure | by reference of another published PhysicalStructureScheme |
| RecordLayoutScheme:<br>• MaintainableType<br>• Description<br>• Label<br>• RecordLayoutSchemeReference<br>• BaseRecordLayout | A separately maintainable scheme, RecordLayoutScheme is basically a listing of individual BaseRecordLayout descriptions included inline or through inclusion by reference of another published RecordLayout Scheme. Note that BaseRecordLayout is an abstract for a substitution group including:<br>• RecordLayout (ASCII fixed or delimited)<br>• DataSet<br>• NCube record layout<br>• Tabular NCube record layout<br>• Inline NCube record layout<br>• Proprietary record layout |

1778

## *8.1   Physical Structure Scheme*

1779

1780  PhysicalStructure is an extension of VersionableType and contains one or more
1781  LogicalProductReferences. Note that a LogicalRecord can contain parts of more
1782  than one LogicalProduct, this element allows a system to identify which
1783  LogicalProduct contents will be referenced from the physical record descriptions.
1784  It also contains one or more GrossRecordStructure descriptions. The remaining
1785  objects within PhysicalStructure are optional and provide default values that
1786  apply to all RecordLayouts referencing a PhysicalRecordSegement described
1787  within this PhysicalStructure unless overridden at a lower level of description.
1788  DDI 3.0 has provided default options a multiple levels and the user needs to
1789  determine how they wish to organize and manage their collection to take
1790  advantage of this feature. For example, if the Format (for example Fixed Format)
1791  is declared at this level, the PhysicalRecordSegments described here cannot be
1792  referenced by a ProprietaryRecordLayout or DataSet. However, it may be
1793  entirely appropriate to define the DefaultDecimalSeparator or DefaultMissingData
1794  here to avoid having to enter this information for each DataItem. Default
1795  declarations are also available all NCube storage descriptions as they can vary
1796  from NCube to NCube. These same elements are found in the PhysicalLocation
1797  of the DataItem description and can be described there. If defaults are used, the
1798  value in the PhysicalLocation will override the default value.
1799

| Element | Usage |
|---|---|
| DefaultDataType | Content is xs:string but preferred use is any W3C datatype |
| DefaultDelimiter | Allowed values are: Empty (default), Tab, Blank, |

| | | AnyString. If a delimiter is used, free field (delimited data) is assumed; binary formats are not allowed. This is the delimiter between data items in a delimited file. |
|---|---|---|
| | DefaultDecimalPositions | This is the default value of implied decimal positions (how many positions to the right of the decimal are included). For example, if the DecimalPosition is 2, then a DataItem value of 8 would be the equivalent of .08. |
| | DefaultDecimalSeparator | There is no default value for this as a "." is common in the US and "," in Europe. |
| | DefaultDigitGroupSeparator | A grouping separator separates groups of digits such as thousands from hundreds. Once again there is no default as "," is common in the US and "." Is common in Europe. |
| | DefaultMissingData | Standardized handling of missing data across the dataset. |

1800
1801
1802 Gross Record Structure provides information on the gross or general physical
1803 structural of a logical record as it appears in a dataset. First is the
1804 LogicalRecordReference which links the physical description to the content
1805 information provided for the LogicalRecord as described in DataRelationship. An
1806 attribute, numberOfPhysicalSegments is set to a default value of one and should
1807 reflect the number of PhysicalRecordSegments defined for the record. You must
1808 describe at least one PhysicalRecordSegment (that is the full record) in order to
1809 provide the link to the RecordLayout.  The minimum description for a logical
1810 record contained in one physical record is an ID (it is an extension of
1811 IdentifiableType), the attribute segmentOrder at its default value of "1" and
1812 attribute hasSegmentKey at the default value of "false".
1813
1814 A logical record that supports multiple segments may be stored as a single
1815 record, multiple segments in a hierarchical file, multiple segments in a data file
1816 per record segment, or a combination of these. Legacy data often has no
1817 structural support for multiple segments but is stored in segments in hierarchical
1818 files where the record order alone defines the segment. If the logical record has
1819 been separated into more than one physical segment the segmentOrder would
1820 be incremented and two optional elements may be used. If there is a segment
1821 Key (a variable that identifies the segement type, the attribute hasSegmentKey is
1822 changed to "true" and the element KeyVariableReference is provided giving the
1823 reference to the Key variable and a value of the variable for the specific segment.
1824
1825 Note that this is a reference is to a single variable so if this is a concatenated
1826 key, a variable whose content is the concatenation of two or more other variables
1827 must be described in the logical record and included in the LogicalRecord. In
1828 addition, a FileNameIdentification string can be provided which indicates how a

1829 segment may be located by its file name structure. For example, in 2000 the U.S.
1830 Census published Summary File 1 in 39 files for the US and for each state. The
1831 file name convention was  xxnnnnn_uf1.zip where xx was either us or the state
1832 two letter postal abbreviation and nnnnn was the segment number.

## *8.2    Record Layout Scheme*

1833
1834 RecordLayoutScheme is an extension of MaintainableType and contains a Label,
1835 Description and means of including one or more record layout structures inline or
1836 by referencing a publish RecordLayoutScheme. The element BaseRecordLayout
1837 is the abstract type for a collection of substitution types to describe the details of
1838 a record layout. This structure allows the development of RecordLayout
1839 descriptions that are specific to the needs of current and future storage systems,
1840 including proprietary systems like statistical software. BaseRecordLayout is an
1841 extension of IdentifiableType and contains a PhysicalStructureReference,
1842 CharacterSet (US ASCII, EBCDIC, UTF-8 etc.), and ArrayBase (1 or 0). If
1843 CharacterSet is unknown, say for a proprietary structure, enter "unspecified".
1844 ArrayBase provides the assumed first position in any array used in describing
1845 this dataset. For example, codebooks historically assume the first character in a
1846 data record is "1", that is, an array base of "1". Unix systems and many
1847 programming languages assume an array base of "0". If the codebook
1848 information uses an array base other than "0" the base level for any array
1849 handled by the system must be declared. If not, all start positions would be offset
1850 by one character. PhysicalStructureReference is a specially structured reference
1851 to a single PhysicalRecordSegment. It contains a reference to the
1852 PhysicalStructure and a separate element, PhysicalRecordSegmentUsed, that
1853 provides the ID of the PhysicalRecordSegment. Note that a BaseRecordLayout
1854 can only link to one PhysicalRecordSegment. If you have a logical record that is
1855 stored in some files as a single record and other files as multiple records, you
1856 must create two different PhysicalStructures to describe the physical stores of
1857 the logical record.
1858

### 8.2.1  RecordLayout

1859
1860 The traditional archive format description of fixed format and delimited files
1861 (similar to earlier DDI versions) is described in RecordLayout as found in
1862 PhysicalDataProduct. All other BaseRecordLayout substitutions will be found in
1863 separate schemas. In addition to the elements inherited from BaseRecordLayout,
1864 RecordLayout allows for a DefaultVariableSchemeReference to limit the need to
1865 repeat this information for each DataItem's VariableReference, an attribute
1866 namesOnFirstRow with a default value of "false", and a list of one or more
1867 DataItem descriptions.
1868
1869 DataItem contains a VariableReference and PhysicalLocation. PhysicalLocation
1870 is used or extended by a number of other BaseRecordLayout substitutions. Its
1871 contents should be familiar to uses of earlier versions of DDI.
1872

| Element | Use |
|---|---|
| StorageFormat | Overrides DefaultFormat. Is a CodeValueType which allows for a controlled vocabulary |
| Delimiter | Overrides DefaultDelimiter |
| StartPosition | Used in fixed format files this is the position of the first character of the data item in the record |
| ArrayPosition | Used in delimited files. This is the array number of the data item. Note that the first item in the array should have a value corresponding to the ArrayBase declared in the RecordLayout. |
| EndPosition | The position of the last character of a data item in a fixed format file |
| Width | The actual width of a data item in a fixed format file or the maximum width in a delimited file. |
| DecimalPostions | Number of decimal places with an implied decimal separator. Another expression of the decimal scalling factor (SAS). Default value is "0" |
| DecimalSeparator | Character used to separate the integer from the faction portion of the number if it is used in the data set. Allowed values include: None (default), Dot, Comma, Other. Overrides DefaultDecimalSeparator |
| DigitGroupSeparator | Character used to separate the sections of an integer (between 100 and 1000 for example) if it is used in the data set. Allowed values include: None (default), Dot, Comma, Other. Overrides Default DigitGroupSeparator |
| LanguageOfData | Use two-character ISO Language Code to indicate the language of the data content. Applicable for text fields. |
| LocaleOfData | Use two-character ISO country code as a supplement to the LanguageOfDataCode |

1873
1874   Note that when describing a fixed format file you must use either EndPosition or
1875   Width in conjunction with StartPostion. You may use both if desired. As you can
1876   see from the number of elements that override defaults that listing a default value
1877   that applies to only half the data items in the records, would considerably reduce
1878   the size and content of the metadata description. In general, it is best to provide

1879    the most common structures at the default level (for example
1880    DataType="Integer") and enter an override value for "String" or "Character" data
1881    items.

## 8.2.2  DataSet

1882

1883    DataSet allows for inline inclusion of data in the metadata file. This is valuable for
1884    small datasets and particularly for small statistical tables. This identifies the
1885    following items:

1886

| Element | Usage |
|---------|-------|
| IdentifyingVariable | References the variable containing the primary key or index value. |
| DefaultVariableScheme | By identifying the variable scheme here, one can enter just the ID for each individual item, reducing repetition. |
| CHOICE: | |
|     RecordSet | Storage structure for a traditional rectangular data structure. Each array of data is for a single record with multiple variables. |
|     ItemSet | Allows data items to be stored in a random order. |
|     VariableSet | Storage structure that is transposed from the traditional record structure. Each array of data is for a single variable with a value for each record in record order. |

1887
1888    The choice between three layouts provides flexibility to store data in whichever
1889    format is optimal for your particular use. A RecordSet requires a list of variables
1890    in order using VariableOrder (a simple list of VariableReference elements listed
1891    in order of appearance in the array. This is followed by Record (repeated for
1892    each record in the dataset) which is the array of values for all variables in the
1893    record. Note that for sparse arrays (arrays with missing values) a separator other
1894    than "blank" must be used in order to express the correct position in the array.
1895    For example:
1896        For Var1=1    Var2=8  Var3=        Var4=10        Var5=
1897        With Blank as delimiter: 1 8  10
1898        With Comma as delimiter: 1,8,,10,
1899
1900    Clearly indicating the missing content for Variable 3 and 5.
1901
1902    ItemSet is a list of ItemValues, each with a VariableReference and a
1903    RecordReference (string containing the value of the IdentifyingVariable), and the
1904    value of the variable.
1905
1906    VariableSet contains a list of VariableItems each containing a VariableReference
1907    and a Value. The assumption is that each record will occupy the same position in

1908   each array and can be identified by the value provided in the VariableItem
1909   referencing IdentifyingVariable. As with RecordSet, datasets with sparse arrays
1910   will need to use a separator other than "blank".

### 1911   8.2.3  NCube Record Layout (Normal)

1912   The first thing to note about NCube physical record layouts is that the
1913   RecordLayout can handle ONLY those dataitems found in NCubes. If your record
1914   has identifying variables fields in addition to the NCube content, you will need to
1915   describe a minimum of 2 PhysicalRecordSegments so that the appropriate
1916   RecordLayout can be used to describe each part. These will later be recombined
1917   in the PhysicalInstance which can contain multiple RecordLayout references
1918   stored as concatenated strings or as hierarchical structures.
1919
1920   The basic NCube structure is for storage structures where one or more NCubes
1921   are stored as records with the data items (cell data) strung out in sequence with
1922   or without a string of non-NCube variables that define the case. Census
1923   aggregate summary files are often structured in this way, where each record is a
1924   geographic area with the data for 100 or more tables strung out in a single logical
1925   record. This type of structure facilitates record subsets as it does not require a
1926   case identification as a dimension of the NCube. This structure can also be used
1927   for cases where all data items are described by the NCube. The NCube record
1928   layout provides a reference to a specific NCube, identification of any attributes
1929   provided for the NCube, and physical location of a data item in a record and
1930   association to the cell coordinates, attribute type, and measure. To identify the
1931   standard content of the NCube such as a count, provide the NCubeReference,
1932   use the Coordinate element to identify the coordinate number as described in the
1933   logical product, and the value of the coordinate or the variable where the value is
1934   being held. In addition you can provide similar location information on attributes
1935   associated with all or parts of the NCube through the use of attribute or
1936   coordinate group. Use coordinate group when the attribute is not associated with
1937   the full NCube.
1938
1939   In addition to the items inherited from BaseRecordLayout, the NCube
1940   RecordLayout is a simple series of NCubeInstance descriptions. Note that an
1941   NCube that is split between 2 or more physical records must be described in 2
1942   separate RecordLayout descriptions. Any cell of an NCube for which there is no
1943   DataItem is assumed to be missing from the PhysicalRecordSegment being
1944   described. It may be missing because it has been declared as being "empty" by
1945   definition or because it is stored elsewhere in another PhysicalRecordSegment.
1946
1947   Each NCubeInstance has a reference to a single NCube description in a
1948   LogicalProduct. It may have an attribute attached at the NCube level for the
1949   NCube as a whole or a specified region of the NCube (as defined in
1950   LogicalRecord). The value of the Attribute may be declared in the metadata or,
1951   using the standard PhysicalLocation structure, reference the location of the
1952   attribute value in the dataset. For example, in many economic tables,

1953 suppression flags are attached to the full table or perhaps to a specific level of
1954 detail for a variable. In these cases a single flag applies to multiple cells in a table
1955 in a consistent fashion for each location. Some locations have suppressed data,
1956 others do not. Attribute at the NCubeInstance level is used only for these
1957 situations, when a single dataitem holds content that applies to multiple cells in
1958 the NCube matrix.
1959
1960 The remainder of the NCubeInstance content is the DataItem along with the
1961 ability to set the standard default values for the NCube as a whole. The DataItem
1962 of an NCube is associated with the NCube by its cell coordinates, the point of
1963 intersection on each of the dimensions of the NCube listed in rank order. The cell
1964 coordinates are provided by the repletion of Dimension for each dimension of the
1965 NCube. Dimension includes an attribute rank whose value corresponds to the
1966 dimension rank of the variable as described in the NCube and either the attribute
1967 value providing the variable value (intersect point) for this cell OR a reference to
1968 a Variable which will contain the value for the intersect point.
1969
1970 Example:
1971
1972 Variable:     Age
1973  1 = Under 18
1974  2 = 18 to 64 years
1975  3 = 65 years and older
1976
1977 Variable:     Sex
1978  1 = Male
1979  2 = Female
1980
1981 Variable:     PoliticalParty
1982 001 = Democrat
1983 002 = Democratic Farmer Labor
1984 003 = Republican
1985 004 = Independent
1986 005 = Green Party
1987 …
1988 560 = Whig
1989
1990 LogicalProduct Description:
1991 NCube1 Age by Sex
1992      Dimension rank="1" VariableReference = "Age"
1993      Dimension rank="2" VariableReference = "Sex"
1994
1995 NCube2 Age by Sex by Political Party Affiliation
1996      Dimension rank="1" VariableReference = "Age"
1997      Dimension rank="2" VariableReference = "Sex"
1998      Dimension rank="3" VariableReference = "PoliticalParty"

1999
2000    For NCube 1 the DataItem for "65 years and older, Female" would be identified
2001    as follow:
2002
2003    <DataItem>
2004            <Dimension rank="1" value="3"/>
2005            <Dimension rank="2" value="2"/>
2006
2007    For NCube 2 there are only as many records for a particular NCube as needed to
2008    capture the Political Parties present for the geographic area in order to avoid a
2009    predominately empty NCube content. For example, Democratic Farmer Labor is
2010    found only in Minnesota. The DataItem for "18 to 64 years, Male, [Any
2011    PoliticalParty]" would be identified as follow:
2012
2013    <DataItem>
2014            <Dimension rank="1" value="2"/>
2015            <Dimension rank="2" value="1"/>
2016            <Dimension rank="2">
2017                    <VariableReference isReference="true">
2018                            </ID>PolticalParty</ID>
2019                    </VariableReference>
2020            </Dimension>
2021
2022    If the value of the DataItem referencing the Variable PolitcalParty is "005" then
2023    the count found in the DataItem described above will be the count of Males 18 to
2024    64 years of age who identified as Green Party, if "002" those who identified with
2025    Democratic Farmer Labor. This structure is used primarily for legacy storage
2026    structures where space was major issue. It is also found in historical files
2027    covering topics that change dramatically in the content of the variable values
2028    over time.
2029
2030    Dimension provides the link between the DataItem being discussed and its
2031    position in the NCube matrix. A DataItem as a cell of an NCube can contain one
2032    or more measures as well as one or more attributes. The attributes described at
2033    this level apply to the specific DataItem only. An example of this is cell level
2034    suppression flags that are stored in the data file. It is identical in structure to that
2035    found at the NCubeInstance level. It provides a reference to the Attribute as
2036    described in the NCube logical product description, information on the physical
2037    storage location using PhysicalLocation or the value of the attribute. In general
2038    the point of having an attribute at this level is that it varies with each instance of
2039    the cell, but the structure allows for declaring a set value in the metadata at this
2040    point.
2041
2042    Measure is also repeatable to allow for multiple measures (count, percent,
2043    cumulative percent, etc.) to be attached to the DataItem. Each Measure contains
2044    a MeasureReference and PhysicalLocation. Note that MeasureReference
2045    extends ReferenceType by providing an attribute arrayOrder (integer with a

2046 default setting of "0" assuming an array base of '0"). Depending on the storage of
2047 the data, multiple measures for a single DataItem can be stored separately or as
2048 an array with a single PhysicalLocation. DDI allows for both storage structures.
2049
2050 If the measures have separate physical location information (different
2051 StartPosition etc.) the element Measure in DataItem is repeated and a single
2052 MearsureReference within Measure is provided. The arrayOrder on
2053 MeasureReference remains at "0".
2054
2055 &lt;DataItem&gt;
2056   &lt;Measure&gt;
2057     &lt;MeasureReference arrayOrder="0"&gt;
2058       &lt;/ID&gt;COUNT&lt;/ID&gt;
2059     &lt;/MeasureReference&gt;
2060     &lt;PhysicalLocation&gt;
2061        .....
2062     &lt;/PhysicalLocation&gt;
2063   &lt;/Measure&gt;
2064   &lt;Measure&gt;
2065     &lt;MeasureReference arrayOrder="0"&gt;
2066       &lt;/ID&gt;PERCENT&lt;/ID&gt;
2067     &lt;/MeasureReference&gt;
2068     &lt;PhysicalLocation&gt;
2069        .....
2070     &lt;/PhysicalLocation&gt;
2071   &lt;/Measure&gt;
2072 &lt;DataItem&gt;
2073
2074 If the measures are stored as an array at a single physical location information
2075 (single StartPosition etc.) the element Measure in DataItem is entered once using
2076 multiple MearsureReferences within Measure and a single PhysicalLocation. The
2077 arrayOrder on MeasureReference would be changed to reflect the position of the
2078 referenced measure in the array.
2079
2080 &lt;DataItem&gt;
2081   &lt;Measure&gt;
2082     &lt;MeasureReference arrayOrder="0"&gt;
2083       &lt;/ID&gt;COUNT&lt;/ID&gt;
2084     &lt;/MeasureReference&gt;
2085     &lt;MeasureReference arrayOrder="1"&gt;
2086       &lt;/ID&gt;PERCENT&lt;/ID&gt;
2087     &lt;/MeasureReference&gt;
2088     &lt;PhysicalLocation&gt;
2089        .....
2090     &lt;/PhysicalLocation&gt;
2091   &lt;/Measure&gt;
2092 &lt;DataItem&gt;
2093

2094     This is the basic structure of NCube storage structures. Tabular and Inline
2095     descriptions will focus on the differences with this basic structure.

2096     **8.2.4 Tabular NCube Record Layout**

2097     A tabular layout is assumed to be a two dimensional layout on a spreadsheet or
2098     print storage that is defined by columns and rows. A spreadsheet can contain
2099     multiple tables on a single sheet and the table may be located anywhere on the
2100     sheet. In addition a "Table" may contain more than a single NCube, either hinged
2101     along a common dimension (see the example in section 7.5) or tightly interlaced
2102     in a specialized layout. The RecordLayout found in
2103     tabular_ncube_recordlayout.xsd can accommodate multiple NCubes found on a
2104     single spreadsheet layout. In addition to the list of NCubeInstances, the tabular
2105     description provides a TopLeftTableAnchor which gives the column and row of
2106     the upper left corner of the table being described. These are expressed as
2107     integers to provide a standard mappable structure for this content.
2108
2109     The NCubeInstance is identical to that of the basic NCube RecordLayout until
2110     one gets to the level of the PhysicalLocation. Tabular does not use the standard
2111     PhysicalLocation, but a specialized extension which adds the elements Column
2112     and RowSequence. Column is the column in which the DataItem will be located.
2113     RowSequence is the Row number within the repeating sequence which holds the
2114     content of the DataItem. For the table Urban/Rural by Age by Nativity by Sex (cell
2115     contents are DataItem coordinate values):
2116

| TABLE 1 | | | Native Born | | Foreign Born | |
|---|---|---|---|---|---|---|
| | | | Male | Female | Male | Female |
| Minnesota | Urban | Under 18 | 1,1,1,1 | 1,1,1,2 | 1,1,2,1 | 1,1,2,2 |
| | | 18 to 64 | 1,2,1,1 | 1,2,1,2 | 1,2,2,1 | 1,2,2,2 |
| | | 65 and over | 1,3,1,1 | 1,3,1,2 | 1,3,2,1 | 1,3,2,2 |
| | Rural | Under 18 | 2,1,1,1 | 2,1,1,2 | 2,1,2,1 | 2,1,2,2 |
| | | 18 to 64 | 2,2,1,1 | 2,2,1,2 | 2,2,2,1 | 2,2,2,2 |
| | | 65 and over | 2,3,1,1 | 2,3,1,2 | 2,3,2,1 | 2,3,2,2 |

2117
2118     This shows a single record sequence so that the Column and RowSequence
2119     values for Rural, 18 to 64, Foreign Born, Male would be Column = 3
2120     RowSequence = 5. Note that the TopLeftTableAnchor would have been stated as
2121     Column=3 and Row=3 (assuming the "TABLE 1" is located at Column=1 and
2122     Row=1 with an array base of 1). If this table is repeated for all states as opposed
2123     to including the states in the NCube structure, the location codes would not
2124     change with the repetitions of each new case.
2125

| TABLE 1 | | | Native Born | | Foreign Born | |
|---|---|---|---|---|---|---|
| | | | Male | Female | Male | Female |
| Minnesota | Urban | Under 18 | 1,1,1,1 | 1,1,1,2 | 1,1,2,1 | 1,1,2,2 |
| | | 18 to 64 | 1,2,1,1 | 1,2,1,2 | 1,2,2,1 | 1,2,2,2 |
| | | 65 and over | 1,3,1,1 | 1,3,1,2 | 1,3,2,1 | 1,3,2,2 |

| | Rural | Under 18 | 2,1,1,1 | 2,1,1,2 | 2,1,2,1 | 2,1,2,2 |
|---|---|---|---|---|---|---|
| | | 18 to 64 | 2,2,1,1 | 2,2,1,2 | 2,2,2,1 | 2,2,2,2 |
| | | 65 and over | 2,3,1,1 | 2,3,1,2 | 2,3,2,1 | 2,3,2,2 |
| Mississippi | Urban | Under 18 | 1,1,1,1 | 1,1,1,2 | 1,1,2,1 | 1,1,2,2 |
| | | 18 to 64 | 1,2,1,1 | 1,2,1,2 | 1,2,2,1 | 1,2,2,2 |
| | | 65 and over | 1,3,1,1 | 1,3,1,2 | 1,3,2,1 | 1,3,2,2 |
| | Rural | Under 18 | 2,1,1,1 | 2,1,1,2 | 2,1,2,1 | 2,1,2,2 |
| | | 18 to 64 | 2,2,1,1 | 2,2,1,2 | 2,2,2,1 | 2,2,2,2 |
| | | 65 and over | 2,3,1,1 | 2,3,1,2 | 2,3,2,1 | 2,3,2,2 |

2126
2127
2128

### 8.2.5  Inline NCube Record Layout

Again, this layout is essentially similar to the basic NCube Record Layout, but in this case rather than a location being specified for a DataItem's attribute and/or measure, the value is provided inline in the instance. This results in the Dataitem gaining an optional attribute xs:lang to provide a language flag at the DataItem level.  There is no PhysicalLocation provided. Attribute retains its optional Value element which becomes required, and Measure replaces PhysicalLocation with a required Value element. Note that Value is of type xs:string to allow for alphanumeric or symbol content. It becomes very important to be sure that all DataItems either use the default data type and other structural defaults, or declares them at the attribute level. Inline NCube RecordLayout is the functional equivalent of DataSet for aggregate data.

### 8.2.6  Proprietary Record Layout (BETA)

This provides a generic structure for describing record layouts for proprietary software, in particular statistical analysis software. In addition to the elements inherited from BaseRecordLayout, ProprietaryRecordLayout includes the following elements.

| Element | Use |
|---|---|
| Software | Standard software identification structure for the software used by this data structure |
| DataItemAddress | Description of how data items are addressed within the file, for example by Variable ID or by Variable Name |
| DefaultNumericDataType | Declares the most common data type used for numeric data using a controlled vocabulary |
| DefaultTextDataType | Declares the most common data type used for text data using a controlled vocabulary |
| DefaultDateTimeDataType | Declares the most common datetime type used for text data using a controlled vocabulary |
| CHOICE: | |

| | |
|---|---|
| CodeDataAsNumeric | Use indicates that variables using CodeRepresentation should normally be treated as numeric data and declares the most common data type using a controlled vocabulary |
| CodeDataAsText | Use indicates that variables using CodeRepresentation should normally be treated as text data and declares the most common data type using a controlled vocabulary |
| ENDCHOICE | |
| ProprietaryInfo | This is a name value pair providing information proprietary to the software package. |
| DataItem | Contains a VariableReference a ProprietaryDataType to override the default, a ProprietaryOutputFormat to designate the desired display or other output format, and ProprietaryInfo at the DataItem level |

2147

## 2148 *7.7   Physical Instance*

2149 Physical Instance has a one-to-one relationship with an actual physical data file.
2150 The single exception is duplicate copies of the same file.  A PhiscalStructure may
2151 have zero (as is the case when using dataset), one, or hundreds of data files
2152 associated with it. Common reasons for this are cases where each record type or
2153 record segment is stored in a separate file or when large datasets result in
2154 subsets of records by spatial, temporal, or topical divisions. For example, the
2155 2000 US Census Summary File 4 has a separate data file for each state by each
2156 record segment by one or a range of characteristic iteration values. This results
2157 in a collection of over one million separate files in order to cover the US. While
2158 this is an extreme case, it is also not uncommon for an archive to acquire a
2159 single dataset such as the Eurobarometer and immediately sort it into a file per
2160 country to facilitate use by their researchers.
2161
2162 Physical Instance is designed to capture these types of divisions as well as
2163 contain the summary and category statistics related to the full dataset and/or
2164 those specific to the particular subset of records.

## 2165 **7.7.1  Top Level Elements**

2166 The Physical Instance has the standard common reusable elements and is the
2167 module that most frequently makes use of Coverage to impose coverage
2168 constraints. The coverage of the Physical Instance will often be a subset of the
2169 study coverage. In the past, coverage constraints of a data file were often only
2170 noted by cryptographic file names. Physical Instance allows for clear evidence of
2171 the subset included in the related data file. Additional elements identify the
2172 Physical Data Product that describes the record(s) contained in the file as well as
2173 the name of the data file itself, its location, and  any additional copies. Attributes
2174 let you indicate which is a master file (as opposed to a backup or other copy),

2175 and note the URI if it is publicly available. In addition to the URI of the data file, a
2176 location and path can be provided as a file may not be available directly due to
2177 access restrictions. (It may, in fact be located on a DVD sitting in a safe, etc.).
2178
2179 In addition to these standard forms of identification, DDI has provided a generic
2180 structure for capturing the "fingerprint" of the data file. This includes a Value for
2181 the fingerprint, the AlgorithmSpecification, and the AlgorithmVersion. The
2182 Fingerprint is repeatable to allow for use of multiple algorithm specifications.

2183 **7.7.2 Gross File Structure**

2184 This section contains information unique to the individual data file and how it was
2185 created.
2186

| ELEMENT | USEAGE |
|---|---|
| PlaceOfProduction | Where the file was produced. |
| ProcessingCheck | Description of any processing checks that were done when the file was made [or subsequent checks]. |
| ProcessingStatus | Many files go through stages in production and the current stage should be noted here. Earlier stages should have been noted in the set of ProcessingCheck entries. |
| CreationSoftware | The name of the software used to create the file including the name, version, description, and date of the software. |
| CaseQuantity | Number of cases in the file. |
| OverallRecordCount | Total number of records in the file (a case may have more than one record). |

2187

2188 **7.7.3 Statistics**

2189 Summary and category statistics for a data file are entered in the physical
2190 instance as the values change when full datasets are subset. Statistics may be
2191 entered inline in the metadata or may exist as a separately described dataset
2192 (common for large complex datasets expressing summary or category statistics
2193 for filtered variables such as each variable by country). The StatisticalDataFile
2194 reference may be to a physical instance that contains the statistics inline or that
2195 represents the data file containing the statistics (noted by an attribute).
2196
2197 Statistics captured inline are organized by variable. The structure provides for
2198 total responses, weight references, handling of missing category, additional
2199 summary statistics, and category statistics. Category statistics can be weighted,
2200 and can be presented as multiple forms of statistics (count, frequency,
2201 cumulative frequency, etc.), and filters. A filter allows you to designate a single
2202 layer cross-tabulation. For example:
2203

| | | |
|---|---|---|
| 2204 | Variable: Sex | COUNT |
| 2205 |     Category: | |
| 2206 |         Male | COUNT |
| 2207 |         Female | COUNT |
| 2208 |      Filter Variable: Country | |
| 2209 |         Germany | |
| 2210 |           Male | COUNT |
| 2211 |           Female | COUNT |
| 2212 |        France | |
| 2213 |           Male | COUNT |
| 2214 |           Female | COUNT |

## 2215 8.0 Group, Resource Package, and Comparison

2216 The schema group.xsd encompasses a number of the features of DDI 3.0 that
2217 allow it to capture the life cycle of data. Group provides an umbrella structure to
2218 pull together two or more studies into a structured series or unstructured group. It
2219 provides basic information on relationships among members of the group that
2220 affect processing decisions based on a required attribute grid (see DDI 3.0
2221 Technical Specification Part I: Overview, Appendix Two) as well as detailed
2222 information on comparable relationships between and among the studies in the
2223 group. In addition, Group contains a specialized structure called
2224 ResourcePackage that allows for publishing maintainable objects (schemes and
2225 schemas) outside of a StudyUnit or Group.
2226

2227 Studies can be grouped for a number of reasons but generally fall into the
2228 category of grouping by design or ad hoc groups. Grouping by design takes
2229 place when studies are either intended to be a series or when a repetition of the
2230 study takes place. The key factor is that the second study in the series is
2231 intended to inherit features of the first study (questions, variables, study design,
2232 universe, etc.) for the purpose of comparability. Group allows you to define which
2233 parts of the major components are shared, where overrides take place, and how
2234 to relate or link data in one study to data in a subsequent survey.
2235

2236 When using inheritance within groups to show comparability – or even just to re-
2237 use metadata – it is important to understand how local overrides work, as this
2238 can impact the way the metadata is grouped. Within each group, all metadata is
2239 inherited down the grouping structure. At any level, it is possible to override any
2240 inherited metadata using the Add, Replace, or Delete attributes which are found
2241 on the IdentifiableType, VersionableType, and MaintainableType structures. To
2242 override an inherited structure, it should have the appropriate ID structure given
2243 for it, and then have the Replace element specified. To delete inherited
2244 metadata, use a similar technique but employ the Delete element. Once replaced
2245 or deleted, it is the modified form of the metadata which is inherited down the
2246 grouping structure.
2247

2248  Note that when referencing metadata that is subject to local overrides, it may be
2249  necessary to specify the exact module being referenced – otherwise, local
2250  deletions and overrides won't be referenced.
2251
2252  Ad hoc groups are collections of studies that have been grouped to meet specific
2253  needs of the archive, data service, or user. These groups do not support
2254  comparison through inheritance as they were developed as separate studies and
2255  grouped later to meet specific needs. Without the use of inheritance,
2256  comparability must be described explicitly using the schema Comparative.
2257  Currently, comparison is enabled for the following complex elements: Universe,
2258  Concept, Question, Category, CodingScheme, and Variable. All comparisons are
2259  pairwise, and with the exception of CodingScheme each comparison notes a
2260  source and target item, the relationship (map) type, and a definition of any
2261  differences when the map type is less than a full equality. CodingScheme
2262  provides options for describing code relationships between two coding schemes.
2263  Options include a human-readable description of the translation process (Source
2264  code 1 through 3 equal Target code 1), a command line for a specified command
2265  language, or the use of the GenerationInstruction from DataCollection. For
2266  example a direct mapping of each source code value to its target value could be
2267  declared.
2268

| **SOURCE** | | **TARGET** |
|---|---|---|
| Code 1 Never Married | recode to | Code 1 Single |
| Code 2 Divorced | recode to | Code 1 Single |
| Code 3 Widowed | recode to | Code 1 Single |
| Code 4 Married | recode to | Code 2 Married |

2269
2270
2271  Alternatiively, GenerationInstruction could be used to identify the Source Variable
2272  and provide the command code If Source >= 1 and <= 3 Target = 1; If Source = 4
2273  Target = 2.
2274
2275  Correspondence does not limit itself to equivalency but captures a text
2276  description of Commonality and Difference as well as a CommonalityTypeCoded,
2277  CommonalityWeight, and a UserDefinedCorrespondenceProperty. The following
2278  table provides some guidance in classifying the level of comparison for Universe,
2279  Concept, Question and Variable as well as the associated Category and Code
2280  Schemes.

 *Data Documentation Initiative*

## Structure of Comparisons

Similar: Denotes a close but not exact relationship; requires description of difference.
Different: Denotes non-equivalent coding scheme; requires coding instructions to create equivalency.

Questions and Variables are assumed to use the categories and coding scheme of the source item.

| Comparison Map | Textual Content of main body | | Category | | Code Scheme | | *ACTION* |
|---|---|---|---|---|---|---|---|
| | *Same* | *Similar* | *Same* | *Similar* | *Same* | *Different* | |
| Universe | X | | ---- | ---- | ---- | ---- | Enter and Flag as Identical |
| | | X | ---- | ---- | ---- | ---- | Flag as similar and note differences in human-readable and optional repeatable machine-actionable string. |
| Concept | X | | ---- | ---- | ---- | ---- | Enter and Flag as Identical. |
| | | X | ---- | ---- | ---- | ---- | Flag as similar and note differences in human-readable and optional repeatable machine-actionable string. |
| Question | X | | X | | X | | Enter and Flag as Identical; include concepts and coding schemes used. |
| | X | | X | | | X | Source and target relationship for question and category scheme must be the same. Reference harmonized coding scheme. |
| | X | | | X | X | | Flag as similar and note category differences in human-readable and optional repeatable machine-actionable string. Assume categories related to codes are those of the source code value. |
| | X | | | X | | X | Flag as similar and note category differences in human-readable and optional repeatable machine-actionable string. Source and target relationship for question and code scheme must be the same. Reference harmonized coding scheme. |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | X | X | | X | | Flag as similar and note text differences in human-readable and optional repeatable machine-actionable string. Assume categories related to codes are those of the source code value. |
| | | X | X | | | X | Flag as similar and note text differences in human-readable and optional repeatable machine-actionable string. Source and target relationship for question and code scheme must be the same. Reference harmonized coding scheme. |
| | | X | | X | X | | Flag as similar and note text and category differences in human-readable and optional repeatable machine-actionable string. Assume categories related to codes are those of the source code value. |
| | | X | | X | | X | Flag as similar and note text and category differences in human-readable and optional repeatable machine-actionable string. Source and target relationship for question and code scheme must be the same. Reference harmonized coding scheme. |
| Variable | X | | X | | X | | Enter and Flag as Identical; include concepts and coding schemes used. |
| | X | | X | | | X | Source and target relationship for variable and code scheme must be the same. Reference harmonized coding scheme. |
| | X | | | X | X | | Flag as similar and note category differences in human-readable and optional repeatable machine-actionable string. Assume categories related to codes are those of the source code value. |
| | X | | | X | | X | Flag as similar and note category differences in human-readable and optional repeatable machine-actionable string. Source and target relationship for variable and code scheme must be the same. Reference harmonized coding scheme. |
| | | X | X | | X | | Flag as similar and note text differences in human-readable and optional repeatable machine-actionable string. Assume categories related to codes are those of the source code value. |

| | | | | | | |
|---|---|---|---|---|---|---|
| | X | X | | | X | Flag as similar and note text differences in human-readable and optional repeatable machine-actionable string. Source and target relationship for variable and code scheme must be the same. Reference harmonized coding scheme. |
| | X | | X | X | | Flag as similar and note text and category differences in human-readable and optional repeatable machine-actionable string. Assume categories related to codes are those of the source code value. |
| | X | | X | | X | Flag as similar and note text and category differences in human-readable and optional repeatable machine-actionable string. Source and target relationship for variable and code scheme must be the same. Reference harmonized coding scheme. |

## 9.0   Step-by-Step Sequence to Create a DDI File for a Simple Instance

These sections should be completed in the following order to ensure that fields containing information that is required (linked) from other elements is available when needed. Note that the list is ordered so that none of the entries requires a section lower down in the list for its creation.

| Sequence No. | Section | Content | Elements/sections referencing items |
|---|---|---|---|
| 1 | **UniverseScheme**<br>All universe items in a structured hierarchy | Universe identification is done by referencing a single universe structure. This may be built up as you go along, but the top level universe should be entered as a container for sub-universes and as the content of the study unit universe. | UniverseReference in studyunit, ControlConstruct, Variable, NCube, UniverseMap. Preexistence of this element also helps in defining coverage and analysis units |

| 2 | **ConceptScheme** A structured list of all concepts of the study | Ideally this should be as detailed as practicable to support ISO 11179. At a minimum there must be at least a single comprehensive concept as ConceptReference is a required element in Question and Variable. | ConceptReference in Question, Variable, and ConceptMap |
|---|---|---|---|
| 3 | **Organization** and **Individual** A listing of all organizations and individuals involved in the life cycle of the study to date | The organization module contains information on organizations, individuals and their structural roles and relationships. This information is housed separately and referenced by a number of element types. | Citation may reference producers, publishers, distributors, etc. Funding information found in studyunit, datacollection, archive, and group reference organizations and individuals. At minimum you will need a listing of the organization or individual who acts as the maintenance agency and for the archive (may be the same). This can be held in a public registry of DDI organizations. |
| 4 | StudyUnit **Citation / Abstract / Purpose** This is a required part of the study unit and represents the fact that a study unit must exist in order to create the specified section | A basic studyunit is the broad description of a simple study. Intellectually, the collection of data and creation of a data file are done within the construct of a study. | The following module cannot exist without a studyunit. Archive must have a study unit to attach to as it describes an archive holding of some study. Group, by definition is two or more study units. |
| 5 | StudyUnit **Coverage** A definition of the topical, temporal, and geographic coverage of the study | Coverage at this level is the top level container describing aspects of coverage. All other instances of coverage are either subsets of this | DataCollection, LogicalProduct, PhysicalDataProduct, PhysicalInstance, Group, and the majority of sub-schemas or sub- |

| | | definition, or in the case of Group, the inclusive combination of the coverage of the study units in the group. Note that if GeographicStructure and GeograhicLocation are used, the Schemes are located in ConceptualComponnents. They may also be published as an external resource. | modules. |
|---|---|---|---|
| **6** | LogicalProduct **CategoryScheme** This contains all categories and their definitions used in the study | Categories are defined once and referenced either directly or through their representational code. | Used to construct Question ResponseDomains, CodeSchemes, and Variables by direct reference or via CodeSchemes. |
| **7** | LogicalProduct **CodeScheme** Connects codes to categories | Organizes categories into structured or unstructured groups and provides the representational code (for example "0 = Male"). For structured CodeSchemes, information on subgroups and relationships is provided. | Question ResponseDomain and Variable Representation |
| **8** | DataCollection **QuestionScheme** Contains question text and response domain information | If Instrument will be included in your instance or if Variable will need to reference the question, the QuestionScheme must be completed. | ControlConstruct, CodingInstructions (ProcessingEvent) (possibly), and Variable (optional) |
| **9** | DataCollection ProcessingEvent **Coding** Explains the process of altering the question response to obtain | This includes general coding instructions, recodes, inclusion of administrative or information from outside of the questionnaire, or | Question (interviewer instructions), ControlConstruct (either), Variable (coding instruction) |

| | | | |
|---|---|---|---|
| | the content of the variable. **InterviewerInstructionScheme** contains all the interviewing instructions including additional descriptive information and visible routing instructions used when completing the questionnaire. | other derived or generated content for the variable. All descriptions are housed in this location and referred to by elements using the process described. All interviewer instructions are held in the scheme and included by reference by the question or the control construct. | |
| 10 | **VariableScheme** Defines the intellectual content and structure of microdata elements and dimensions for aggregate data matrixes (NCubes) | The variable scheme defines individual variables. They are used to construct NCubes and provide a single source of intellectual content regardless of the storage structure. Variables can be grouped in several ways including those contained within a single type of record. | NCube, VariableGroup, DataRelationships, DataItem, SummaryStatistic, CategoryStatistic |
| 11 | LogicalProduct **DataRelationships** Defines the identification of each record type, information needed to identify a specific case, and links between record types | The structural information provided here is a map to assist in record identification, selection, and linking. This is the intellectual map of the types of structures and linkages supported by the variable content of the record. Actual physical layout is described separately, but use the information provided here, simply adding the specifics of how the physical file dealt with the intellectual structure. | PhysicalStructure. Required link to the LogicalRecord |
| 12 | **NCubeScheme** Describes aggregate data | The NCube is defined by its universe and dimensions. Each | DataItem |

| | | cell of the NCube is identified by its matrix coordinate pattern. | |
|---|---|---|---|
| **13** | **PhysicalStructure/ PhysicalRecordSegment** Identifies each physical storage record type and how it uses the intellectual links between record types as described in DataRelationships | A data set is made up one or more physical record types stored in one or more physical data file structures. | RecordLayout contains a required link between the RecordLayout and the PhysicalRecordSegment. |
| **14** | **RecordLayout** desctibes the physical layout of data items within a record OR provides the data inline | The DataItems describe the physical location of the variable or NCube cell being described on a specified physical record type. The PhysicalInstance may contain one or multiple physical record types and must identify which of these types described in PhysicalData that it contains. | PhysicalInstance contains a required link to one or more RecordLayouts. |

The following diagram identifies sections of DDI in the order they need to be completed. The gray boxes are major steps and include one or more sections of DDI elements. The solid yellow boxes represent element sets that will be required for references or other use later in creating the full instance. The boxes with yellow diagonal bars are required if you provide references (such as a reference from a variable to a question) or have this feature (NCubes). Arrows indicate references from one element set to another. Dotted line arrows indicate that the inclusion of this information is by choice. If you do not have the element set you do not need the references. References from STEP 8 are dependent upon what features you choose to include; however, at this point the required material is available for reference.

**<ddi>** *Data Documentation Initiative*

**STEP 1**
- Universe Scheme
- Concept Scheme
- Organization Scheme
- StudyUnit Citation/ Abstract/ Purpose
- StudyUnit Coverage

**STEP 2**
- Category Scheme
- Coding Scheme

**STEP 3 optional**
- Question Scheme
- Control Construct Scheme
- Coding and Interviewer Instruction Scheme

**STEP 4**
- Variable Scheme
- Data Relationships
- NCube
- Remaining Logical Product items

**STEP 5**
- Physical Structure
- Record Layout

**STEP 7**
- Physical Instance

**STEP 8**
- Archive / Group / etc.