

Using the SAS® Clinical Standards Toolkit for define.xml creation

Lex Jansen, SAS Institute Inc., Cary, NC

ABSTRACT

When submitting clinical study data in electronic format to the FDA, not only information from trials has to be submitted, but also information to help understand the data. Part of this information is a data definition file, which is the metadata describing the format and content of the submitted data sets. When submitting data in the CDISC SDTM format it is required to submit the data definition file in the Case Report Tabulation Data Definition Specification (CRT-DDS, also known as define.xml) format as prepared by the CDISC define.xml team.

This workshop will provide an introduction to the structure and content of the define.xml file. The SAS® Clinical Standards Toolkit will then be used to create the define.xml file.

INTRODUCTION

The FDA issued the Final Guidance on Electronic Submissions using the eCTD specifications in April 2006. The latest revision of this guidance was published in June 2008 [1]

Technical specifications associated with this guidance are provided as stand-alone documents. Among these are Study Data Specifications [2] that provide further guidance for submitting animal and human study data in electronic format when providing electronic submissions to the FDA. Study data includes information from trials submitted to the agency for evaluation and information to understand the data (data definition). The study data includes both raw and derived data.

As of January 1, 2008, sponsors submitting data electronically to the FDA are required to follow the new eCTD guidance. The previous guidance, originally issued in 1999, has been withdrawn as of the same date.

The new guidance differs from the 1999 guidance in one significant aspect: The application table of contents is no longer submitted as a PDF file, but is submitted as XML (eXtensible Markup Language). This means that the electronic submissions will now be XML based.

The current version of the Study Data Specifications contains specifications for the Data Tabulation data sets of human drug product clinical studies and provides a reference to the Study Data Tabulation Model (SDTM) [3][4][5][6] developed by the Submission Data Standard (SDS) working group of the Clinical Data Interchange Standard Consortium (CDISC).

Further, the Study Data Specifications document gives a reference to the Case Report Tabulation Data Definition Specification (CRT-DDS or define.xml) developed by the CDISC define.xml Team [7].

DATA DEFINITION TABLES: define.xml

Released for implementation in February 2005, the Case Report Tabulation Data Definition Specification (CRT-DDS or define.xml) Version 1.0 [7] specifies the standard for providing Case Report Tabulations Data Definitions in an XML format for submission to regulatory authorities (e.g., FDA). The XML schema used to define the expected structure for these XML files is based on an extension to the CDISC Operational Data Model (ODM). The current version of the CRT-DDS (version 1.0) is based on version 1.2.1 of the CDISC ODM [8].

The Data Definition Document provides a list of the data sets included in the submission along with a detailed description of the contents of each data set. To increase the level of automation and improve the efficiency of the Regulatory Review process, the define.xml file can be used to provide the Data Definition Document in a machine-readable format.

XML 101

In this section we present a short introduction to XML.

BASIC SYNTAX

The Extensible Markup Language (XML) [9] is a general-purpose markup language. It is classified as an extensible language because it allows users to define their own elements. Its primary purpose is to facilitate the sharing of

structured data across different information systems. XML is a language that is hierarchical, text-based and describes data in terms of markup tags. A good introductory guide to XML can be found in the reference section [10].

Every XML document starts with the **XML declaration**. This is a small collection of details that prepares an XML processor for working with the document.

syntax	XML declaration	example
	<code><?xml param1 param2 ... ?></code>	<code><?xml version="1.0" encoding="ISO-8859-1" ?></code>

Elements are the basic building blocks of XML, dividing a document into a hierarchy of regions. Some elements are containers, holding text or (child) elements. Others are empty and serve as place markers. Every XML file should have exactly 1 root element that contains all other elements.

syntax	Container Element	example
	<code>< name attribute1 attribute2 ... > content </ name ></code>	<code><def:leaf ID="Location.DM" xlink:href="dm.xpt"> <def:title>dm.xpt</def:title> </def:leaf></code>

An **empty element** is similar, but contains no content or end tag.

syntax	Empty Element	example
	<code>< name attribute1 attribute2 ... /></code>	<code><ItemRef ItemOID="STUDYID" OrderNumber="1" Mandatory="Yes" Role="IDENTIFIER" RoleCodeListOID="RoleCodeList" /></code>

In the element starting tag there can be information about the element in the form of an **attribute**. Attributes define properties of elements. They associate a name with a value, which is a string of character data enclosed in quotes. There is no limit to how many attributes an element can have, as long as no two attributes have the same name.

syntax	Attributes	example
	<code>name = " value "</code>	<code>ItemOID="STUDYID" OrderNumber="1" Mandatory="Yes" Role="IDENTIFIER" RoleCodeListOID="RoleCodeList"</code>

Namespaces are a mechanism by which element and attribute names can be assigned to groups. They are most often used when combining different vocabularies in the same document. If each vocabulary is given a namespace then the ambiguity between identically named elements or attributes can be resolved.

syntax	Namespaces	example
	<code>xmlns: name = " URI "</code>	<code>xmlns="http://www.cdisc.org/ns/odm/v1.2" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xmlns:xlink="http://www.w3.org/1999/xlink" xmlns:def="http://www.cdisc.org/ns/def/v1.0"</code>

Comments in an XML document are not interpreted by the XML processor:

syntax	Comments	example
	<code><!-- comment text --></code>	<code><!-- File: define.xml --> <!-- Date: 02/02/2005 --> <!-- Author: Clinical Data Interchange Standards --> <!-- Consortium (CDISC) define.xml team --></code>

A **CDATA** (character data) section tells the XML parser that this section of the document contains no markup and should be treated as regular text. CDATA sections can be useful for large regions of text that contain a lot of 'forbidden' XML characters. They should be used with care though, since it may be hard to use any elements or attributes inside the marked region.

syntax	CDATA	example
	<pre><![CDATA unparsed character data]]></pre>	<pre><TranslatedText xml:lang="en"><![CDATA[Classifies subjects based on Hy's Law - 2 variations - For each variation, subject is classified as normal or abnormal at baseline and as normal or abnormal based on most extreme value during treatment, Safety population]]></TranslatedText></pre>

Processing instructions are meant to provide information to a specific XML processor, but may not be relevant to others. It is a container for data that is targeted toward a specific XML processor. The processing instruction looks like the XML declaration, but is targeted at a specific XML processor. The XML declaration can be viewed as a processing instruction for all XML processors.

syntax	Processing Instructions	example
	<pre><? target data ?></pre>	<pre><?xml-stylesheet type="text/xsl" href="define1-0-0.xsl" ?></pre>

Now that we have explained the most important building blocks of XML, we can look at a more complete example in Figure 1.

Figure 1: Excerpt of an XML file (define.xml)

```
<?xml version="1.0" encoding="ISO-8859-1" ?>
<?xml-stylesheet type="text/xsl" href="define1-0-0.xsl"?>
<ODM xmlns="http://www.cdisc.org/ns/odm/v1.2"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xmlns:xlink="http://www.w3.org/1999/xlink"
  xsi:schemaLocation="http://www.cdisc.org/ns/def/v1.0
  http://www.cdisc.org/ns/odm/v1.2 define1-0-0.xsd"
  FileOID="Studydisc01"
  ODMVersion="1.2" FileType="Snapshot"
  CreationDateTime="2007-04-09T12:24:09">
  - <Study OID="cdisc01">
  - <GlobalVariables>
    <StudyName>CDISC01</StudyName>
    <StudyDescription>CDISC01 Test Study.</StudyDescription>
    <ProtocolName>CDISC01</ProtocolName>
  </GlobalVariables>
  <MetaDataVersion OID="CDISC.SDTM.3.1.1"
    Name="Study CDISC01, Data Definitions"
    Description="Study CDISC01, Data Definitions"
    def:DefineVersion="1.0.0"
    def:StandardName="CDISC SDTM"
    def:StandardVersion="3.1.1">
  ...
  <ItemGroupDef OID="MH"
    Name="MH" Repeating="Yes" IsReferenceData="No"
    Purpose="Tabulation" def:Label="Medical History"
    def:Structure="One record per medical history event per subject"
    def:DomainKeys="STUDYID, USUBJID, HHCAT, HHTERM, MHSTDTCT"
    def:Class="EVENTS"
    def:ArchiveLocationID="Location.AE">
  <ItemRef ItemOID="STUDYID" OrderNumber="1" Mandatory="Yes"
    Role="IDENTIFIER" RoleCodeListOID="RoleCodeList" />
  ...
  <def:leaf ID="Location.MH" xlink:href="mh.xpt">
    <def:title>mh.xpt</def:title>
  </def:leaf>
</ItemGroupDef>
```

The first line is the XML declaration.

The second line contains a processing instruction to include the XSL style sheet to display the XML file in a human readable form in a browser.

<ODM> is the root element.

<Study>, <GlobalVariables>, <MetaDataVersion>, <StudyName>, <ItemGroupDef>, etc are elements.

<Study> is a child element of the <ODM> element.

FileOid, ODMVersion, etc are attributes of the <ODM> element.

The 4 lines starting with "xmlns:" are namespace declarations. This excerpt contains no comments or CDATA sections.

WELL-FORMED AND VALIDATED

An XML file is said to be *well-formed* if it conforms to the rules of XML syntax. At a very basic level this means:

- A single element, called the root element, contains all other elements in the document (in the define.xml the root element is <ODM>)
- Elements should be properly opened and closed

- Elements do not overlap, e.g. are properly nested
- Attributes are properly quoted
- The document does not contain illegal characters.
For example, the “<” character has special meaning because it opens a tag. So, if this character is part of the content, it should be substituted as “<”

A conforming XML parser is not allowed to process an XML document that is not well-formed.

An **XML schema** is a description of a type of XML document, typically expressed in terms of constraints on the structure and content of documents of that type. This description goes above and beyond the basic syntax constraints imposed by well-formedness of an XML document [11].

The schema defines the allowed elements and attributes, order of the elements, overall structure, etc...

A schema might also describe that the content of a certain element is only valid if it conforms to the ISO 8601 date and time specification. An XML document is *valid*, if it conforms to a specific XML schema.

The following example illustrates the difference between **well-formed** and **validated**. The XML document in Figure 2 is a well-formed XML document, but is obviously not valid with respect to the schema that defines a valid define.xml file.

Figure 2: A well-formed XML document, but not a valid define.xml document¹

```
<?xml version="1.0" encoding="ISO-8859-1" ?>
<heartofrockandsoul>
  <single genre="country">
    <rank>435</rank>
    <artist>George Jones</artist>
    <title>He Stopped Lovin' Her Today</title>
    <producer>Billy Sherrill</producer>
    <writer> Bobby Braddock & Curly Putnam</writer>
    <label>Epic 50867</label>
    <year>1980</year>
    <billboard>Did not make pop charts</billboard>
  </single>
  <single genre="country">
    <rank>928</rank>
    <artist>Jerry Lee Lewis</artist>
    <title> One Has My Name (the Other Has My Heart)</title>
    <producer>Jerry Kennedy</producer>
    <writer> Eddie Dean, Dearest Dean & Hal Blair</writer>
    <label>Smash 2224</label>
    <year>1969</year>
    <billboard> Did not make pop charts</billboard>
  </single>
</heartofrockandsoul>
```

XPATH AND XSL

In order to be able to understand the XML structure of a define.xml file, we need to briefly discuss a few more XML related standards and concepts.

The XML Path Language (**XPath**) gives XML developers a tool for navigating the structure of an XML document. We can simply demonstrate this by 2 examples from Figure 2.

- The XPath location path `/heartofrockandsoul/single/artist` returns all artist elements: “George Jones” and “Jerry Lee Lewis”.
- The XPath location path `/heartofrockandsoul/single/@genre` returns the “country” attribute.

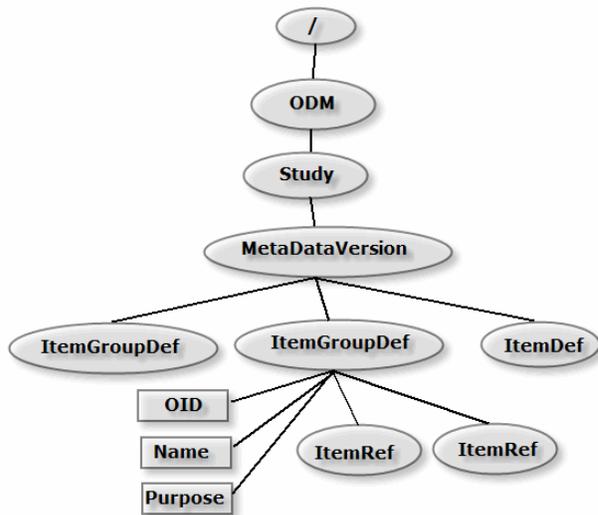
XPath sees an XML document as a tree with branches called nodes. The nodes in this tree have family relationships: parent-child, ancestor-descendant, sibling, and so forth.

¹ This particular document may not be a valid define.xml document, but it is certainly appropriate for a PharmaSUG in Nashville, home of the Grand Ole Opry.

Conceptually, the XPath data model describes an XML document as having seven possible node types:

- Root
- Element
- Attribute
- Text
- Namespace
- Comment
- Processing instruction

Figure 3: define.xml as a tree



In Figure 3 we see part of a define.xml as a tree with nodes. We only show element and attribute nodes. **ODM** is the root element and is the *parent* of the **Study** element. The **ItemGroupDef** and **ItemDef** elements are *siblings*. The **OID**, **Name** and **Purpose** attributes and the **ItemRef** elements are *children* of The **ItemGroupDef** element. **Study**, **MetaDataVersion**, **ItemGroupDef**, **ItemDef** and **ItemRef** are all *descendants* of the ODM element. The **ODM** element is an *ancestor* of the **ItemRef** elements, or the **ItemDefs** elements.

XPath allows us to access any node from any other node simply by knowing the relationship between the two.

At the foundation of the XPath language is the ability to use *location paths* to refer to a node or nodeset. A node is a piece of the XML document (such as an element, an attribute, or text content).

A location path identifies a set of nodes in a document. This set may be empty, may contain a single node, or may contain several nodes. These can be element nodes, attribute nodes, namespace nodes, text nodes, comment nodes, processing-instruction nodes, root nodes, or any combination of these. A location path is built out of successive *location steps*. Each location step is evaluated relative to a particular node in the document called the *context node* [15].

There are two kinds of location paths: relative location paths and absolute location paths.

A *relative location* path consists of a sequence of location steps separated by a forward slash (/). Each step selects a node or nodeset relative to the current node. Then, each node in that set is used as the current node for the following step, and so on.

An *absolute location* path starts with a forward slash, optionally followed by a relative location path.

Examples:

- odm:Decode/odm:TranslatedText/text()
- /odm:ODM/odm:Study/odm:MetaDataVersion/ItemGroupDef

The SAS XML Mapper application also uses XPath location paths to tell the SAS XML engine where to find content from the XML elements and attribute content to create rows and columns in SAS data sets.

Figure 4 shows part of an XMLMAP. It specifies that the data set **ItemDef** has a variable named **Name**, whose content is defined by the XPath specification **/LIBRARY/ItemDefs/Name**.

Figure 4: XMLMAP example

```

<TABLE name="ItemDefs">
  <TABLE-PATH syntax="XPath">/LIBRARY/ItemDefs</TABLE-PATH>
  <TABLE-DESCRIPTION></TABLE-DESCRIPTION>

  <COLUMN name="OID" > [6 lines]
  <COLUMN name="Name">
    <PATH syntax="XPath">/LIBRARY/ItemDefs/Name</PATH>
    <TYPE>character</TYPE>
    <DATATYPE>character</DATATYPE>
    <DESCRIPTION></DESCRIPTION>
    <LENGTH>128</LENGTH>
  </COLUMN>
  <COLUMN name="DataType" > [6 lines]
  <COLUMN name="Length" > [6 lines]
  <COLUMN name="SignificantDigits" > [6 lines]
  <COLUMN name="SASFieldName" > [6 lines]
  <COLUMN name="SDSVarName" > [6 lines]
  <COLUMN name="Origin" > [6 lines]
  <COLUMN name="Comment" > [6 lines]
  <COLUMN name="CodeListRef" > [6 lines]
  <COLUMN name="Label" > [6 lines]
  <COLUMN name="DisplayFormat" > [6 lines]
  <COLUMN name="ComputationMethodOID" > [6 lines]
  <COLUMN name="FK_MetaDataVersion" > [6 lines]

</TABLE>

```

The XPath standard is part of the XSL family of standards.

XSL is the Extensible Style sheet Language [12], one of the most complicated – and most useful – parts of XML. While XML itself is intended to define the structure of a document, it does not contain any information on how it is to be displayed. In order to do this we need a language, XSL, to describe the format of a document, ready for use in a display application (computer screen, cell phone screen, paper).

XSL is actually a family of transformation languages which allows one to describe how files encoded in the XML standard are to be formatted or transformed.

The following three languages can be distinguished:

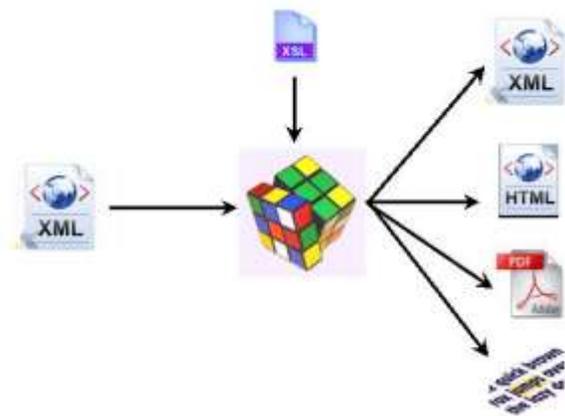
- XSL Transformations (XSLT): an XML language for transforming XML documents into other XML documents, into HTML documents, or into any other text based documents (like a SAS program), or even a PDF file.
- XSL Formatting Objects (XSL-FO): an XML language for specifying the visual formatting of an XML document
- The XML Path Language (XPath): a non-XML language used by XSLT, and also available for use in non-XSLT contexts, for addressing the parts of an XML document.

An XSL stylesheet can be used to transform the define.xml to HTML in such a way that it can be viewed in a browser.

Figure 5: Define.xml displayed in a browser using an XSL style sheet

Variable	Label	Type	Controlled Terminology	Origin	Role	Comment
STUDYID	Study Identifier	text		CRF Page 3	IDENTIFIER	
DOMAIN	Domain Abbreviation	text		DERIVED	IDENTIFIER	
USUBID	Unique Subject Identifier	text		SPONSOR DEFINED	IDENTIFIER	Concatenation of STUDYID.SUBJID
SUBJID	Subject Identifier for the Study	text		CRF Page 3	TOPIC	

Figure 6: Transforming with XSLT



THE STRUCTURE OF THE DEFINE.XML

The previous section explained the building blocks of XML. This section will specifically describe the structure of a define.xml file that conforms to the Case Report Tabulation Data Definition Specification (CRT-DDS or define.xml) version 1.0 as developed by the CDISC define.xml Team [7].

SCHEMA STRUCTURE

The CRT-DDS (define.xml) standard is based on the CDISC operational model (ODM). The ODM is defined by an XML schema that allows extension [8]. This extension mechanism has been implemented by expressing the ODM schema using two files:

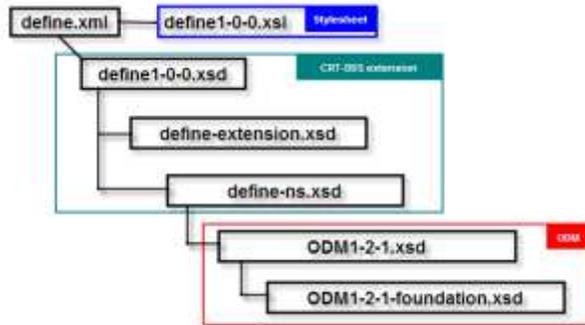
- a *foundation* XML Schema file (ODM1-2-1-foundation.xsd), which defines the elements, attributes and structure of the base ODM schema.
- an *application* XML Schema file (ODM1-2-1.xsd) which imports the foundation XML Schema and other schema definitions needed by ODM, such as the core W3C XML schema (xml.xsd) and the XML schema that defines the W3C XML Signature standard (xmldsig-core-schema.xsd).

To create the define.xml extension three files have been provided:

- A *namespace* XML schema (define-ns.xsd) that defines the extension namespace and any new elements and attributes.
- An *extension* XML Schema file (define-extension.xsd) defines the location of the extensions within the ODM.
- An *application extension* XML Schema file (define1-0-0.xsd) that will import the extension XML Schema file and, in turn, any files imported in the ODM root schema.

In November 2009 CDISC published a XML Schema Validation for Define.xml white paper [13]. This white paper, created by the CDISC XML Technology team, provides guidance on validating define.xml documents against the define.xml XML schemas and proposes practices and tools to improve define.xml schema validation.

Figure 7: The define.xml (CRT-DDS) and the associated style sheet and schemas



DEFINE.XML BUILDING BLOCKS

In this paragraph we will show the building blocks of the define.xml. As mentioned before, the CRT-DDS standard (define.xml) is an extension of ODM. The ODM model defines many elements and attributes that are optional.

We will only mention required elements and attributes, or elements and attributes that were found in the example define.xml file.

Figure 8 shows a high-level overview of a define.xml file.

Figure 8: The building blocks of the define.xml (CRT-DDS)



We will now describe the different parts of the define.xml file.

1 The **header section** of the define.xml is important for the XML processor.

The first line identifies the file as an XML document and specifies the XML version ("1.0") and the encoding of the document ("ISO-8859-1"). The second line includes a reference to the style sheet ("define1-0-0.xsl") that can be used by an XSL processor to render the document. In this case the style sheet can be used by the XSL processor in a web

browser to render the XML file as HTML for display. As mentioned before, the HTML that gets created by the browser depends on the particular browser application.

② Following the header section is the **ODM** root element. All other elements in the define.xml file will be contained within the ODM element. The ODM element contains attributes that define the namespaces and the location of the schema that specifies the XML document.

Other required ODM attributes are displayed in Table 1.

Table 1: ODM required attributes

Attribute	Description
FileType	Type of transmission. For define.xml the only valid value is "Snapshot". This means that the document contains only the current state of the data and metadata it describes, and no transactional history.
FileOID	A unique identifier for this file. FileOIDs should be universally unique if at all possible.
CreationDateTime	Date and Time when the define.xml file was last modified (ISO 8601 datetime).
ODMVersion	The version of the ODM standard used. According to the ODM schema, this is an optional attribute. However, a missing ODMVersion should be interpreted as "1.1". Documents based on ODM 1.2 should have ODMVersion="1.2".

③ **Study** is the first element contained in the ODM root element. The Study element collects static structural information about an individual study. It has one attribute "OID", which is the unique identifier of the Study.

The Study element has two child-elements in the define.xml:

- GlobalVariables - General summary information about the Study.
- MetaDataVersion - Domain and variable definitions within the submission.

④ **GlobalVariables** is a required child element of the **Study** element and contains three required child elements:

- StudyName - Short external name for the study.
- StudyDescription - Free-text description of the study.
- ProtocolName - The sponsor's internal name for the protocol.

⑤ The **MetaDataVersion** is a child element of the **Study** element and contains the domain and variable definitions included within a submission. Table 2 lists the MetaDataVersion attributes that are part of the define.xml file. The attributes with a prefix of "def:" are extensions to the ODM schema.

Table 2: MetaDataVersion required attributes

Attribute	Description
OID	The unique identifier of the MetaDataVersion
Name	Name of the MetaDataVersion
Description	Further description of the data definition document
def:DefineVersion	The schema version used for the define.xml
def:StandardName	Short name of the MetaDataVersion (e.g. CDISC SDTM)
def:StandardVersion	The version of an external standard to which the data conforms (e.g. 3.1.1)

⑥ Table 3 lists the MetaDataVersion child elements. The ItemGroupDef and ItemDef elements are required.

Table 3: MetaDataVersion child elements

Element	Description
def:AnnotatedCRF	This element can be used to reference an Annotated Case Report Form (CRF), a PDF file that maps the data collection fields used to collect subject data to the corresponding variables or discrete variable values contained within the data sets
def:SupplementalDoc	This element can be used to reference a PDF file with supplemental data definition information. A reason for this document can be the need for further explanations or descriptions of variables contained within the data sets.
def:leaf	This element has to be present if either the def:AnnotatedCRF and/or the def:SupplementalDoc are provided. The def:leaf element will then contain the actual location of the PDF file relative to the define.xml
def:ComputationMethod	This element is available for every unique computational algorithm that is referenced by variable metadata or variable value-level metadata. It contains the method name and the computation rule in a kind of pseudo code.
def:ValueListDef	This element provides additional value-level metadata for certain variables that are part of a normalized data structure. For example, the Vital Signs domain has a measurement parameter (VSTESTCD) that stores the name of a measurement (height, weight, systolic blood pressure, ...). A corresponding value list will then describe each unique value of the measurement ("HEIGHT", "WEIGHT", "SYSBP", ...)
ItemGroupDef	For every data set in the submission there is an ItemGroupDef element describing data set metadata (label, structure, keys, variables, location of transport file, ...)
ItemDef	For every variable referenced in either def:ValueListDef or ItemGroupDef there will be an ItemDef element. This element will describe properties like name, label, length, computation method, range or code list restrictions, and several other properties.
CodeList	The CodeList element defines a discrete set of permitted values for an item. The definition can be an explicit list of values or a reference to an externally defined code list.

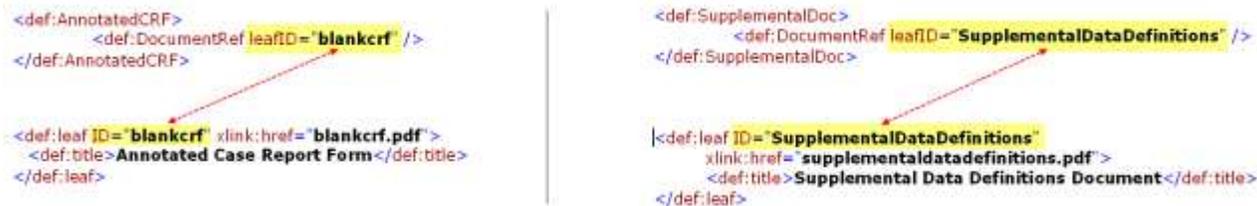
METADATAVERSION ELEMENT DETAILS

In this paragraph we will dive deeper into the sub-elements of the **MetaDataVersion** element. We will see how the different elements relate to each other.

def:AnnotatedCRF, def:SupplementalDoc and def:leaf

As mentioned before, these elements can be used to reference PDF files. Figure 9 shows the relation between these elements. The **def:DocumentRef/@leafID** must match the corresponding **def:leaf/@ID**.

Figure 9: Referencing external PDF files



ItemGroupDef

For every data set in the submission there will be an ItemGroupDef element with domain-level metadata. Table 4 lists the ItemGroupDef attributes.

Table 4: *ItemGroupDef* attributes

Attribute	Status	Description
OID	Required	The unique identifier of the domain.
Name	Required	The file name of the data set or data domain name (e.g., "DM" for Demographics)
Repeating	Required	"Yes" for domains with more than one record per subject whereas, "No" for domains with 1 record per subject.
IsReferenceData	Optional	"Yes" for domains that contain reference data only, no subject-level data. "No" indicates subject-level data. Absence of this attribute indicates subject-level data.
SASDatasetName ²	Optional	Name of SAS data set.
Purpose	Optional	Purpose for the data set (e.g., "Tabulation", "Analysis").
def:Label	Required	Brief description of the data domain.
def:Structure	Optional	Data domain structure (e.g. "One record per subject per visit").
def:DomainKeys	Optional	Comma-separated text string that contains the data set variables that uniquely identify a data record.
def:Class	Optional	General class of the domain as defined in the SDTM model (e.g., "Events", "Interventions", "Findings", "Special Purpose", "Trial Design", ...)
def:ArchiveLocationID	Required	Reference to the def:leaf element that contains a link to the location of the data set.

ItemGroupDef elements have:

- One ItemRef child element per variable in the data set
- Exactly one def:leaf element that contains the XLink information that is referenced by the def:ArchiveLocationID attribute.

The relation between the **ItemGroupDef/@def:ArchiveLocationID** and the **def:leaf/@ID** is illustrated in Figure 10.

Figure 10: *Referencing external data sets*

```

<ItemGroupDef OID="MH"
  Name="MH"
  Repeating="Yes"
  IsReferenceData="No"
  SASDatasetName="MH"
  Purpose="Tabulation"
  def:Label="Medical History"
  def:Structure="One record per medical history event per subject"
  def:DomainKeys="STUDYID, USUBJID, MHCAT, MHTERM, MHSTDTCT"
  def:Class="EVENTS"
  def:ArchiveLocationID="Location.MH">
  <ItemRef ... />
  ...
  <ItemRef ... />
  <def:leaf ID="Location.MH" xlink:href="mh.xpt">
    <def:title>mh.xpt</def:title>
  </def:leaf>
</ItemGroupDef>

```

ItemRef and ItemDef

For every variable in a data set in the submission there will be an ItemRef and an ItemDef element with variable-level metadata. Table 5 lists the ItemRef attributes. Table 6 lists ItemDef attributes. ItemDef elements can have an optional associated CodeListRef or def:ValueListRef child element.

² SASDatasetName is an optional attribute that is part of the ODM foundation. This means that it is a valid attribute in the define.xml. This attribute is not mentioned in the define.xml specification.

def:ComputationMethod

The def:ComputationMethod element has one attribute (OID) and contains the details about computational algorithms used to derive or impute variable values.

Figure 11: Example of a ComputationMethod

```
<def:ComputationMethod OID="COMPMETHOD.QTCB">  
  QTcB = QT interval / square root of (60 / heart rate)  
</def:ComputationMethod>
```

Table 5: ItemRef attributes

Attribute	Status	Description
ItemOID	Required	The unique identifier of the variable.
OrderNumber	Optional	Provide an ordering on the ItemRefs (within the containing ItemGroupDef) for use whenever a list of ItemRefs is presented to a user.
Mandatory	Required	Indicates that the clinical data for an instance of the containing item group would be incomplete without an instance of this type of item. Variables that have an SDTM "Core" attribute given as "Required" should have Mandatory="Yes" in the define.xml. SDTM variables that have a "Core" attribute as "Expected" or "Permissible" should have Mandatory="No" in the define.xml.
Role	Optional	Variable classification (e.g., "Identifier", "Topic", "Timing Variable", ...). Values are interpreted as codes from CodeList as specified by @RoleCodeListOID
RoleCodeListOID	Optional	The identifier of the corresponding Role Code List, which defines the full set roles from which the Role attribute values are to be taken.

Table 6: ItemDef attributes

Attribute	Status	Description
OID	Required	The unique identifier of the variable.
Name	Required	Variable name.
DataType	Required	Type of the data (e.g., text, integer, float, ...)
Length	Conditional	The variable length.
SignificantDigits	Conditional	The number of decimal digits.
SASFieldName ³	Optional	
Origin	Optional	Indicator of the origin of the variable (e.g., CRF page number, derived, or a reference to other variable(s)).
Comment	Required	Further information regarding variable definition, usage, etc.
Def:Label	Required	Variable label.
Def:DisplayFormat	Optional	Display format for numeric variables (e.g., 8.2)
Def:ComputationMethodOID	Optional	Unique identifier of the corresponding computation method.

CodeList and def:ValueListDef

Many variables in a submission have a discrete list of valid values or controlled terms associated with them. The CodeList element defines the controlled terminology. For variables whose values are restricted to a list of values, the corresponding ItemDef element has to include a CodeListRef element that references a CodeList element. The CodeListOID attribute of the CodeListRef element should match the OID attribute of the CodeList element.

³ SASFieldName is optional an attribute that is part of the ODM foundation. This means that it is a valid attribute in the define.xml. This attribute is not mentioned in the define.xml specification.

The `def:ValueListDef` element provides additional value-level metadata for certain variables that are part of a normalized data structure. For example, the Vital Signs domain has a measurement parameter (`VSTESTCD`) that stores the name of a measurement (height, weight, systolic blood pressure, ...). A corresponding value list will then describe each unique value of the measurement (`"HEIGHT"`, `"WEIGHT"`, `"SYSBP"`, ...). For these variables, the corresponding `ItemDef` element has to include a `def:ValueListRef` element that references a `def:ValueListDef` element. The **`def:ValueListOID`** attribute of the **`def:ValueListRef`** element should match the **`OID`** attribute of the **`def:ValueListDef`** element.

Figure 12 illustrates several concepts that are related to `ItemRef`, `ItemDef`, `def:ComputationMethod`, `CodeList` and `def:ValueListDef` elements.

The black numbers on the right side represent the source of a relation and a red numbers on the left side represent the target of a relation. In every relation there has to be a correspondence between the identifier of the source and the target.

For example: In relation 4: source identifier = **`ItemRef/@ItemOID`**, target identifier = **`ItemDef/@OID`** and the both have the value `"VS.VSTESTCD"`;

Figure 12: Relations within a define.xml



DESIGNING A RELATIONAL DATA STRUCTURE FOR THE DEFINE.XML

In the previous paragraphs we have seen the structure of the define.xml file. Also, we have seen the relations between the different parts within the define.xml file. Since the define.xml file does not have a 2-dimensional data structure, it is not a trivial task to translate the define.xml to a number of SAS data set with rows and columns. SAS has defined a relational data model that represents the define.xml.

The highly structured nature of CDISC CRT-DDS data requires that any mapping to a relational format include a large number of data sets, with foreign key relationships to help preserve the intended non-relational object structure. In the SAS Clinical Standards Toolkit, foreign key relationships are enforced when validating the CDISC CRT-DDS data sets. Appendix 1 and Appendix 2 show the 39 SAS data sets that represent the CRT-DDS data model.

One could use SAS/Base to construct the information that needs to go into the 39 tables, however, it is important to understand what the relation is between the CRT-DDS tables and what information is needed in the tables. Some of the information in these 39 tables can be easily derived from the SDTM datasets that need to be described by the define.xml (dataset label, dataset name, variable names and labels, ...). Other information can not automatically be created and needs to be created through a process.

THE SAS CLINICAL STANDARDS TOOLKIT

The SAS Clinical Standards Toolkit is a technology that exists as a set of SAS macros and metadata SAS datasets. The macros that are part of the toolkit are used by SAS® Clinical Data Integration, which can be used to populate metadata datasets and perform macro calls through a user interface.

VALIDATING THE DEFINE.XML

Earlier we talked about two different ways to verify whether XML files are coded correctly:

Well-formed The XML code must be syntactically correct.

Valid If the XML file has an associated XML Schema, the elements must appear in the defined structure and the content of the individual elements must conform to the declared data types specified in the schema.

There are numerous tools available that can check well-formedness and validate an XML file against an XML schema. As mentioned earlier in this paper, the CDISC XML Technologies Team has published a XML Schema Validation for Define.xml white paper, that provides guidance on validating define.xml documents against the define.xml XML schemas and proposes practices and tools to improve define.xml schema validation [13].

To further ensure the quality of the define.xml, more validation needs to be performed:

- Validate against the rules as specified in the Case Report Tabulation Data Definition Specification (define.xml) version 1.0 [7]. An example of this is the rule that an **ItemRef** element must be uniquely defined within an **ItemGroupDef**, through the uniqueness of the **ItemOID** attribute of the **ItemRef**.
- Check the consistency between the define.xml content and the SAS transport files. For example, an **ItemGroupDef** element must contain an **ItemRef** element for each variable included in the corresponding SAS transport file.
- Check the consistency between the define.xml and the SDTM metadata. An example of this is the "Mandatory" attribute of the **ItemRef** element, which needs to be consistent with the SDTM Core attribute (Required, Expected, Permissible).

Once we have the define.xml converted into relational SAS data sets, it will be easy to perform these validation checks using the SAS Clinical Standards Toolkit.

USING THE SAS CLINICAL STANDARDS TOOLKIT

Before the macros in the SAS Clinical Standards Toolkit can be used there are a few tasks that need to be performed.:

- Setup of global macro variables.

This can be done either directly through **%let** statements, or by loading properties (used to create macro variables) with the **cst_setStandardProperties** macro:

```
%cst_setStandardProperties(_cstStandard=CST-FRAMEWORK,_cstSubType=initialize);  
%cst_setStandardProperties(_cstStandard=CDISC-CRTDDS,  
_cstStandardVersion=1.0,_cstSubType=initialize);
```

- Process input and output has to be defined by creating a process **SASReferences** data set. This data set is then communicated with the **cstutil_processsetup** macro.

An example of a SASReferences dataset for the process that creates a define.xml file is presented in Figure 13.

Figure 13: SASReferences data set

VIEWTABLE: Work.Sasreferences									
	standard	standardversion	type	subtype	SASref	reftype	path	order	memname
1	CST-FRAMEWORK	1.2	messages		messages	libref		1	
2	CDISC-CRTDDS	1.0	messages		crmsg	libref		2	
3	CDISC-CRTDDS	1.0	autocall		auto1	fileref		1	
4	CDISC-CRTDDS	1.0	control	reference	control	libref	C:\Users\frjans\AppData\Local\Te Temporary Files_TD6716		. sasreferences
5	CDISC-CRTDDS	1.0	results	results	results	libref	C:\HOW\jansen\cdisc-crtdds-1.0/r		. write_results.sas7bdat
6	CDISC-CRTDDS	1.0	sourcedata		srodata	libref	C:\HOW\jansen\cdisc-crtdds-1.0/		.
7	CDISC-CRTDDS	1.0	externalxml	xml	extxml	fileref	C:\HOW\jansen\cdisc-crtdds-1.0/s		. define.xml
8	CDISC-CRTDDS	1.0	referencexml	stylesheet	xslt01	fileref			.
9	CDISC-CRTDDS	1.0	properties	initialize	inprop	fileref		1	initialize.properties

The following toolkit macros are typically used in the process to create a define.xml:

- **crtdds_sdtm311todefined10**
This macro creates the 39 CRT-DDS data sets by using SDTM table, column, and study metadata.
- **crtdds_validate**
This macro validates the CRT-DDS data sets to make sure that they are in compliance with business rules associated with the CRT-DDS standard
- **crtdds_write**
This macro creates the XML rendition of the CRT-DDS dataset, i.e. the define.xml
- **crtdds_xmlvalidate**
This macro perform XML schema validation of the define.xml file

GENERATING ADDITIONAL METADATA FOR THE CRT-DDS FILE

Some of the information that needs to be represented in the define.xml file cannot be automatically generated and needs to be provided by the user. The user will need to add metadata directly to the CRT-DDS datasets, and then generate the define.xml with the **crtdds_write** macro to include this additional metadata.

An example of this is the so-called value level metadata, which is needed for SDTM Findings domains like Subject Characteristics (CS) domain. The SCTEST variable contains values like "Subject Initials" and "Race". We need to define metadata based on the value of SCTEST:

Value Level Metadata							
Source Variable	Value	Label	Type	Controlled Terms or Format	Length	Role	Comment
SCTEST	Initials	Initials	text		3		Initials
SCTEST	Race	Race	text		20		Race, Other

To be able to represent this information we will need to add data to several datasets, as can be seen in Figure 14.

The following datasets are involved:

- ItemDefs** contains the column metadata
- ValueLists** contains id of value lists
- ValueListItemRefs** contains id of each item in a value list – associates each value list item to a row in the ItemDefs table
- ItemValueListRefs** associates each value list item to a row in the ItemDefs table

Figure 14: CRT-DDS Datasets related to value level metadata

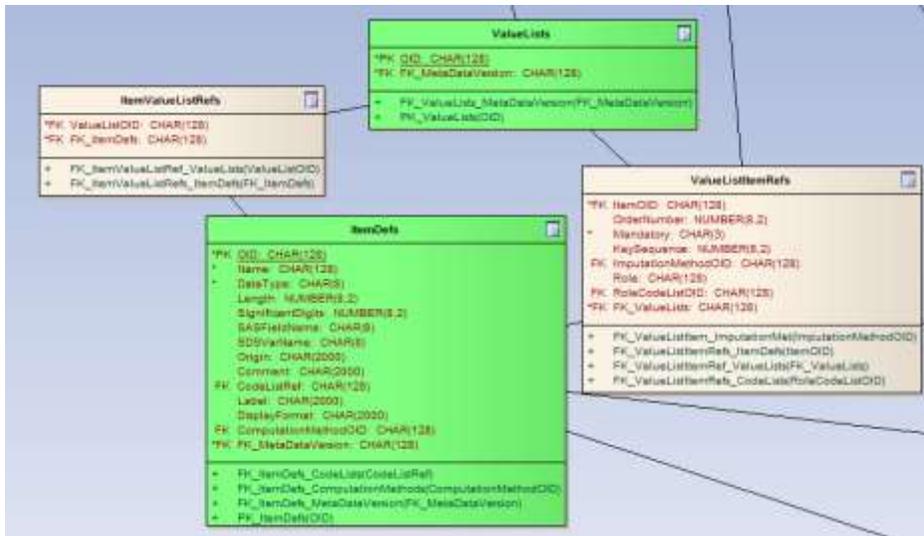


Figure 15: Sample PROC SQL code to provide valuelevel metadata

```

*Lookup OID for the SCTEST column in ItemDefs;
proc sql noprint;
  select OID into :srccol from srcupd.ItemDefs
  where name='SCTEST';
quit;

*add record for a new valuelist SCTESTVALS;
proc sql ;
  insert into srcupd.ValueLists
  set OID= "SCTESTVALS",
  FK_MetaDataVersion = "&mdv";
*add record associating the value list SCTESTVALS to the OID for SCTEST ItemDefs record;
insert into srcupd.ItemValueListRefs
  set ValueListOID= "SCTESTVALS",
  FK_ItemDefs = "&srccol";
*add records to the ItemDefs dataset for each value in the SCTESTVAL value list;
insert into srcupd.ItemDefs
  set OID= "VAL001", Name = "Initials",
  DataType = "text", Length = 3,
  SASFieldName = "INITIALS",
  comment = "Initials",
  label="Initials", FK_MetaDataVersion = "&mdv"
set OID= "VAL002", Name = "Race",
  DataType = "text", Length = 20,
  SASFieldName = "RACEOTH",
  comment = "Race, Other",
  label="Race", FK_MetaDataVersion = "&mdv";
*add records associating the value list SCTESTVALS to rows in the ItemDefs dataset;
insert into srcupd.ValueListItemRefs
  set ItemOID= "VAL001",
  OrderNumber=1, Mandatory="Yes",
  KeySequence=1, FK_ValueLists = "SCTESTVALS"
  set ItemOID= "VAL002",
  OrderNumber=2, Mandatory="Yes",
  KeySequence=2, FK_ValueLists = "SCTESTVALS";
quit;

```

CONCLUSION

The SAS Clinical Standards Toolkit provides metadata and macros to work with the CRT-DDS dataset representation of the define.xml. The Toolkit is able to transform metadata into a define.xml that can be validated for both internal consistency and consistency with the xml CRT-DDS schema.

REFERENCES

1. U.S. Department of Health and Human Services, Food and Drug Administration, Center for Drug Evaluation and Research (CDER), Center for Biologics Evaluation and Research (CBER)". "Final Guidance for Industry: Providing Regulatory Submissions in Electronic Format--Human Pharmaceutical Applications and Related Submissions Using the eCTD Specifications". Revision 2, June 2008. (<http://www.fda.gov/downloads/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/UCM072349.pdf>)
2. U.S. Department of Health and Human Services Food and Drug Administration Center for Drug Evaluation and Research (CDER). Study Data Specifications, Version 1.5.1, January 2010 (<http://www.fda.gov/downloads/Drugs/DevelopmentApprovalProcess/FormsSubmissionRequirements/ElectronicSubmissions/UCM199759.pdf>)
3. CDISC Study Data Tabulation Model, Version 1.1, April 28, 2005 (<http://www.cdisc.org/content1605>)
4. CDISC Study Data Tabulation Model Implementation Guide: Human Clinical Trials, Version 3.1.1, August 26, 2005 (<http://www.cdisc.org/content1605>)
5. CDISC Study Data Tabulation Model, Version 1.2, November 12, 2008 (<http://www.cdisc.org/sdtm>)
6. CDISC Study Data Tabulation Model Implementation Guide: Human Clinical Trials, Version 3.1.2, November 12, 2008 (<http://www.cdisc.org/sdtm>)
7. Case Report Tabulation Data Definition Specification (define.xml), Version 1.0, February 9, 2005 (<http://www.cdisc.org/define-xml>)
8. CDISC Operational Data Model (ODM), Version 1.2.1, January, 2005 (<http://www.cdisc.org/odm>)
9. Extensible Markup Language (XML) 1.0, Fourth Edition, August 16, 2006 (<http://www.w3.org/TR/2006/REC-xml-20060816>)
10. Eric T. Ray, 2003, [Learning XML](#), *Creating Self-Describing Data*. 2nd Edition, (O'Reilly and Associates)
11. Eric van der Vlist, 2002, [XML Schema](#), *The W3C's Object-Oriented Descriptions for XML* (O'Reilly and Associates)
12. Doug Tidwell, 2001, [XSLT](#), *Mastering XML Transformations*. (O'Reilly and Associates)
13. XML Schema Validation for Define.xml, Version 1.0, November 30, 2009 (<http://www.cdisc.org/define-xml>)
14. SAS Institute Inc. 2010. SAS® Clinical Standards Toolkit 1.3: User's Guide . Cary, NC: SAS Institute Inc.
15. Gene Lightfoot (2010). Implementing, Managing, and Validating a Clinical Standard Using SAS Clinical Standards Toolkit 1.3. Proceedings of the 6th Pharmaceutical Users Software Exchange (PhUSE 2010, Berlin, Germany) (<http://www.lexjansen.com/phuse>)
16. Julie Maddox, Mark Lambrecht (2010). define.xml – Tips and Techniques for creating CRT-DDS. Proceedings of the 6th Pharmaceutical Users Software Exchange (PhUSE 2010, Berlin, Germany) (<http://www.lexjansen.com/phuse>)

CONTACT INFORMATION

Lex Jansen,
SAS Institute Inc.
919-531-9860
Email: <mailto:lex.jansen@sas.com>

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration. Other brand and product names are trademarks of their respective companies.

APPENDIX 1

SAS Data Sets representation of CRT-DDS (define.xml)

DefineDocument	contains information specifically about the define.xml file, including description and file creation date/time
Study	contains descriptive information about each study included in the define.xml file
MetaDataVersion	contains information about the data standards used in the study.
AnnotatedCRF	contains document references to the annotated case report form
SupplementalDocs	contains document references to the supplemental documentation
MDVLeaf	contains the href link for each of the documents listed in the AnnotatedCRF and SupplementalDocs tables.
MDVLeafTitles	contains a descriptive title for each MDVLeaf item.
ProtocolEventRefs	contains a list of study events and the order they occur
StudyEventDefs	contains descriptive information on the study events
StudyEventFormRefs	associates a particular form to a specific study event
FormDefs	contains a list of forms for the study
FormDefItemGroupRefs	associates a table with a form definition
FormDefArchLayouts	contains document references for the form layouts
Presentation	contains how information on the study is presented
ItemGroupDefs	contains metadata for each of the tables included in the study
ItemGroupAliases	contains context and alias names for the table
ItemGroupLeaf	contains href location to the specified table
ItemGroupLeafTitles	contains descriptive title to table location
ItemGroupDefItemRefs	associates the columns to the study tables
ItemDefs	contains the column metadata
ItemQuestionTranslatedText	contains translated text of questions
ItemQuestionExternal	contains the name and version of external question dictionaries
ItemMURefs	associates a column to a measurement unit
ItemRangeChecks	contains information on range checks for data in a column
ItemRangeCheckValues	contains information on range checks for data in a column
RCErrrorTranslatedText	error message to be used when a range check is violated
ItemRole	contains the role of the column
ItemAliases	contains context and alias names for the column
MeasurementUnites	contains metadata about measurement units
MUTranslatedText	contains translated text of measurement units
CodeLists	contains name, type and SAS Format name of code lists
ExternalCodeLists	contains name and versions of external code lists
CodeListItems	contains a list of coded items for each code list
CLItemDecodeTranslatedText	contains the full text for a coded item in the specified language
ImputationMethods	contains description of imputation methods
ComputationMethods	contains description of computation methods
ValueLists	contains id of value lists
ValueListItemRefs	contains id of each item in a value list – associates each value list item to a row in the ItemDefs table
ItemValueListRefs	associates each value list item to a row in the ItemDefs table

APPENDIX 2

Data model for CRT-DDS (define.xml)

