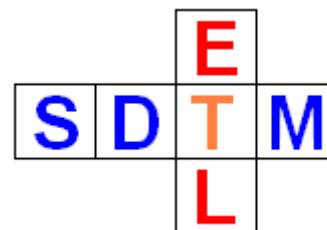# SDTM-ETL 4.1 User Manual and Tutorial

Author: Jozef Aerts, XML4Pharma

Last update: 2022-09-24

## Generating datasets in the CDISC Dataset-JSON format - the submission format of the future

SAS Transport 5 ("XPT format") is a very old format, one generation after the punch card, and thus completely outdated. Still the regulatory authorities such as FDA, PMDA and NMPA require the use of the SAS Transport 5 format for regulatory submissions. As other industries have already since a long time given up the use of SAS Transport (5 or 8) format, the Clinical Research industry is the only industry that is still using this outdated format.

A few of the most important disadvantages of SAS Transport are:
- only US-ASCII is supported.
  This is especially important for submissions to PMDA and NMPA (no support for Japanese and Chinese characters)
- 8-, 40- and 200-character (bytes) limitations for variable names, variable labels and variable values.
- not vendor-neutral
- not very suitable for modern RESTful web services

And in the scope of CDISC standards
- 8-character (bytes) limitation for variables values that might be used in transposal such as --TESTCD, QNAM, ...
- 40-characters (bytes) limitation for variables values that might be used as labels in transposal such as --TEST, QLABEL, ...

Some have proposed the use of SAS Transport _8_ as an solution for the character length limitations but this essentially is a fake solution, as it doesn't solve the ASCII-only problems[1] and vendor-neutrality issues.

Already many years ago, CDISC developed Dataset-XML as a modern alternative for SAS Transport. The FDA then conducted a pilot with a number of pharma companies and with the support of CDISC. Essentially however, the pilot approach was not very good, essentially testing features that were already proven many times before. For example, the feature that Dataset-XML fully supports Unicode (which supports all written languages in the world) was not taken into account at all.
The FDA report (unfortunately not available anymore) mentioned that the FDA would expect problems with file sizes, as in some cases, Dataset-XML files can be considerably larger in size than SAS-XPT files. We think however that this was just a cheap excuse to not have to change anything.

Recently, CDISC developed the Dataset-JSON format, which is similar to Dataset-XML, but then in ... JSON. Major advantages of Dataset-JSON are:
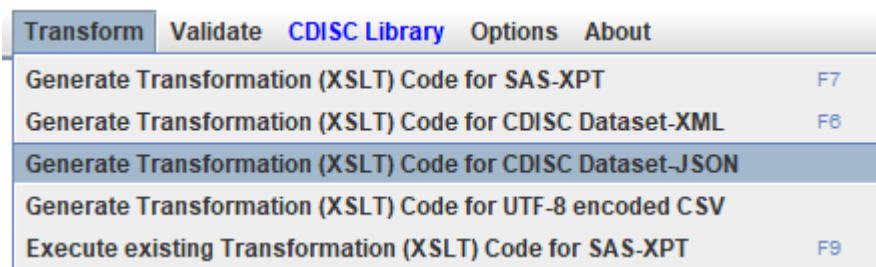
---

[1] Using SAS Transport 8, it is possible to have non-ASCII characters using some "tricks", but when inspecting such files, it is impossible to find out whether, and which of these "tricks" have been used.

- All advantages of Dataset-XML: no character limitations, Unicode "out of the box"
- Much smaller file sizes than SAS Transport 5 (even by a factor of 10)
- For basic consumption, no use of the define.xml is necessary
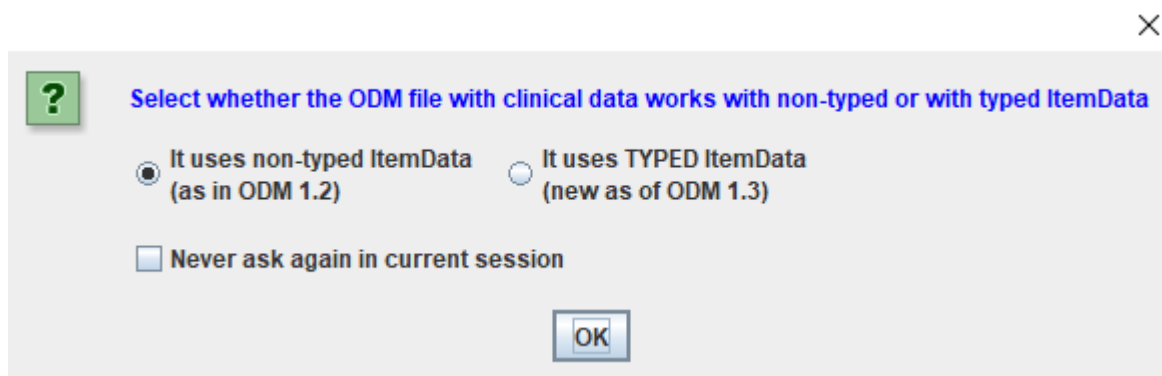- Extremely suitable for use with RESTful web services and APIs

The latter would be extremely worthful once regulatory authorities and the pharma industry come to a "neutral zone", very probably in the cloud, where submissions can essentially be started once the first data point has been collected, and which is governed by RESTful web services, i.e. having the ability to work without the concept of "files".

SDTM-ETL now fully supports the generation of SDTM and SEND files in the new CDISC Dataset-JSON format, including the visualization with the free and open-source "Smart Submission Dataset Viewer".

In order to generate datasets in the new Dataset-JSON, instead of using the menu "Transform - Generate Transformation (XSLT) Code for SAS-XPT", the menu "Transform - Generate Transformation (XSLT) Code for CDISC Dataset-JSON":



After a few moments, the following dialog then shows up:
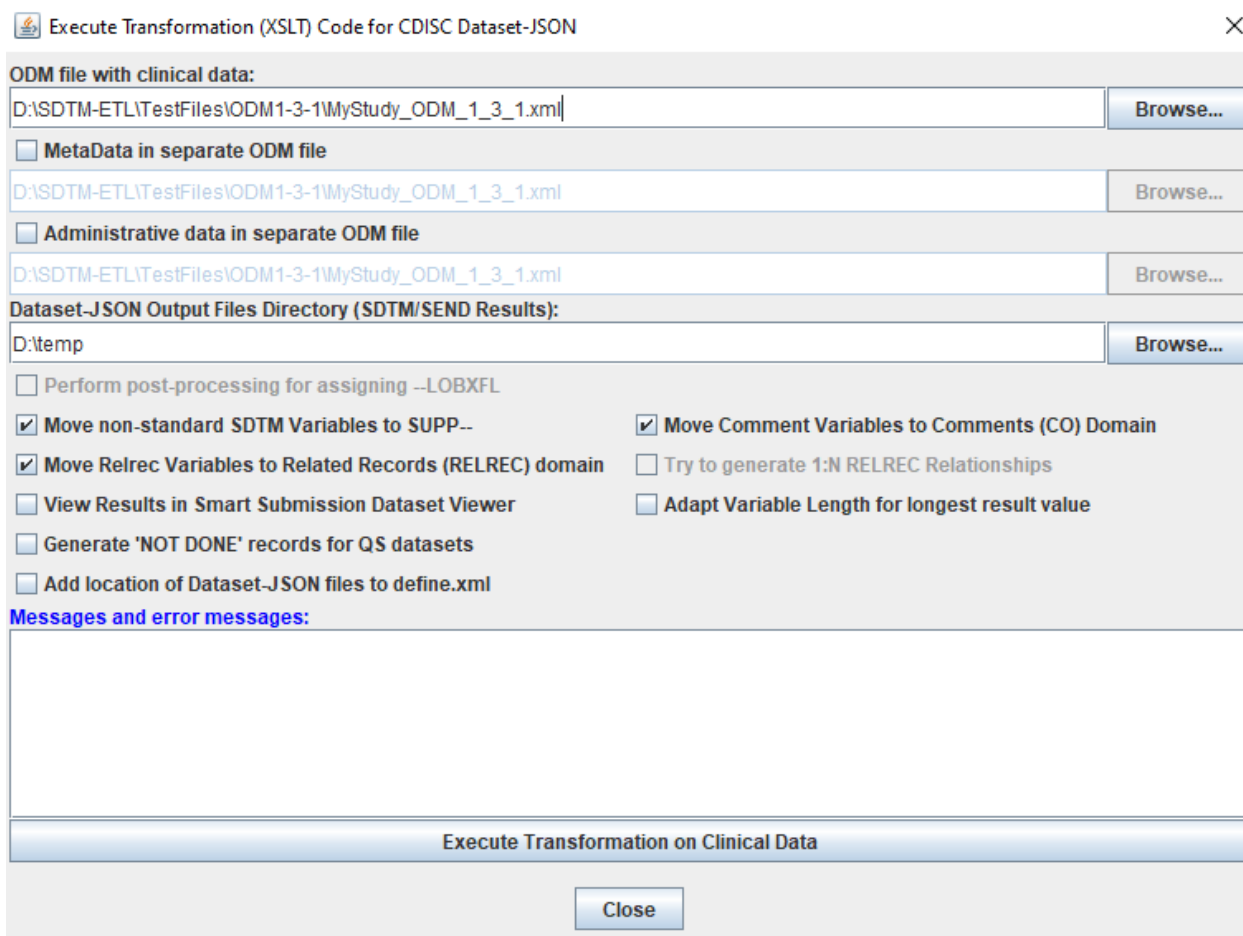


Allowing to choose between the two "tastes" of ODM files with clinical data, and then followed by a new window:

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform"
3 xmlns:sdm="http://www.cdisc.org/ns/studydesign/v1.0"
4 xmlns:math="http://www.w3.org/2005/xpath-functions/math"
5 xmlns:xdt="http://www.w3.org/2005/02/xpath-datatypes"
6 xmlns:odm="http://www.cdisc.org/ns/odm/v1.3"
7 xmlns:xs="http://www.w3.org/2001/XMLSchema"
8 xmlns:sdtm-etl="http://www.xml4pharma.com/SDTM-ETL/ns"
9 xmlns:rws="http://www.xml4pharma.com/SDTM-ETL/RWS/ns"
10 xmlns:fn="http://www.xml4pharma.com/SDTM-ETL/functions/local"
11 xmlns:data="http://www.cdisc.org/ns/Dataset-XML/v1.0"
12 xmlns:OpenClinica="http://www.openclinica.org/ns/odm_ext_v130/v3.1"
13 xmlns:def="http://www.cdisc.org/ns/def/v.2.1"
14 version="2.0">
15
16
17 <xsl:output method="xml" encoding="UTF-8" indent="yes" />
18 <xsl:variable name="SINGLE_QUOTE">&apos;</xsl:variable><xsl:variable name="DOUBLE_QUOTE">&quot;</xsl:variable><xsl:variable name
19 <xsl:param name="LOINC2SDTMLB_CSVFILELOCATION">D:/eclipse-java-2018-09-win32-x86_64/eclipse/workspace/SDTM-ETL_4_1/CDISC
20 <xsl:param name="LOINC2SDTMLB_XMLFILELOCATION">D:/eclipse-java-2018-09-win32-x86_64/eclipse/workspace/SDTM-ETL_4_1/CDISC
21 <xsl:output name="sds-xml-format" method="xml" indent="yes"/><xsl:template match="/"><xsl:apply-templates/></xsl:template>
22 <!-- Template for the top ODM element -->
23
24 <xsl:template match="odm:ODM">
25 <!-- XSLT generated from the SDTM-ETL scripting language - domain instance = MyStudy:GLOBAL --><!-- XSLT generated from the SDTM-ETL
26 <xsl:element name="ODM" xmlns="http://www.cdisc.org/ns/odm/v1.3" xmlns:data="http://www.cdisc.org/ns/Dataset-XML/v1.0" xmlns:xlink="http
27
28 <xsl:attribute name="ODMVersion">1.3.2</xsl:attribute>
29 <xsl:attribute name="Description">SDTM data generated by the SDTM-ETL tool</xsl:attribute>
30 <xsl:attribute name="FileType">Snapshot</xsl:attribute>
31 <xsl:attribute name="FileOID">MyStudy:AE</xsl:attribute>
32 <xsl:attribute name="data:DatasetXMLVersion">1.0.0</xsl:attribute>
33 <!-- Add an instruction that automatically creates a datetime stamp when the stylesheet is executed -->
34 <xsl:attribute name="CreationDateTime"><xsl:value-of select="current-dateTime()"/></xsl:attribute>
35 <!-- Add a ClinicalData element -->
```

The generated transformation code can then be saved to file for later execution, e.g. when batch generation of SDTM/SEND datasets are planned (e.g. using a job-scheduling mechanism), or executed when clicking the "Execute Transformation (XSLT) Code". The latter then leads to:

The window looks very similar to the one when generating SAS-XPT files with some major differences:

- There is no checkbox "Split records > 200 characters to SUPP-- records"
  The reason is that, as a modern format, Dataset-JSON does not have any character length limitations, so that no splitting is at all necessary
- The checkbox "View Result SDTM tables" is replaced by a checkbox "View Results in Smart Submission Dataset Viewer.
  If this checkbox is checked, after the datasets are generated, the "Smart Submission Dataset Viewer" will be started, allowing to inspect the datasets, and much much more (as it is "smart").
- The field for providing the location where the generated datasets need to be written to (a folder) is provided higher up.
- There checkbox "Adapt Variable Length for longest result value" is still present. In the case of SAS-XPT it is necessary for making the XPT files as compact as possible (SAS-XPT is an extremely inefficient format). In the case of Dataset-JSON, it is still there to set the "Length" attribute on the variable definition in the define.xml, as this is still required by the Define-XML standard. In Dataset-JSON, adding the maximal length is optional, but in our case, it is still added, in order to provide backward compatibility for some tools.

We add some information and select some of the choices, e.g.:

asking to put the results in the folder "D:\MyStudy_Dataset-JSON", and to open the "Smart Submission Dataset Viewer" once the generated files are ready.

Then clicking "Execute Transformation on Clinical Data" starts the transformation.
As we have checked the checkbox "Adapt Variable Length", first a window is displayed, allowing us to select for which variables we want to have the maximal length adapted:

During mapping development, one will often choose to not adapt variable lengths, but it can be of importance when generating the final datasets.

In our case, remark the "highest" length being 349 for PEREASND. In the case of SAS-XPT, this would have led to two additional records in an artificially generated SUPPPE.xpt file. In the case of Dataset-JSON, everything can stay nicely together.

If we check all the checkboxes and then click "OK", some messages first pass by:
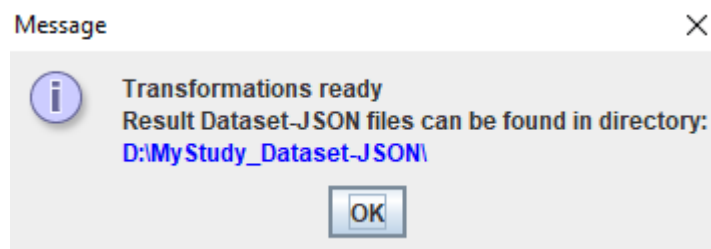


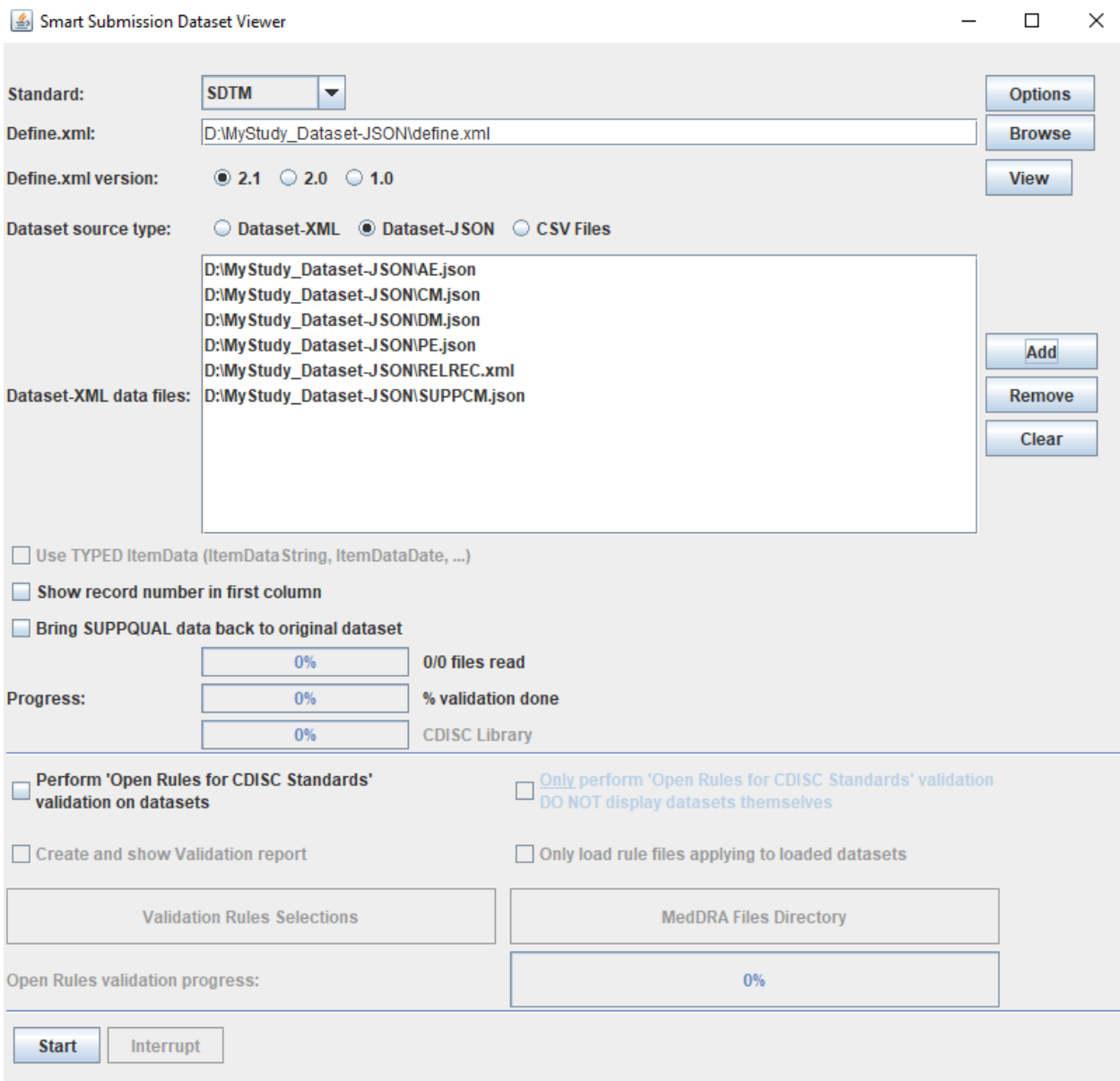and some intermediate information and question messages, like:

which is mostly important in the case of Define-XML 2.1.
Followed by:



and:



The "Smart Submission Dataset Viewer" then opens with:

Listing the generated datasets and the location of the define.xml.

One can then switch on many of the "smart" features of the "Smart Submission Dataset Viewer" by clicking the "Options" button. During generation of the mappings and testing, this will however usually not be necessary.

What is often interesting however is to check the checkbox "Bring SUPPQUAL data back to original dataset", which will generate a "view" where SUPPQUAL records with NSVs (Non-Standard-Variables) are recombined with the original records.

Then clicking "Start" starts the generation and then leads to a set of tables like:



Some of the features that are always interesting to use are:

- when selecting a row in a non-DM table, CTRL-D will select and show the corresponding DM record. Using CTRL-D one can than "toggle" between both the datasets
- when selecting a row in a SUPPxx dataset, CTRL-S will select and show the corresponding source record(s). Also here, one can the toggle between the two datasets by the use of CTRL-B.

If we had checked the checkbox "Bring SUPPQUAL data back to original dataset", the view for the CM dataset would become:



The Dataset-JSON files themselves are very compact (file sizes considerably smaller than for SAS-XPT). For example, for DM.json:



In future, once a "neutral zone" between sponsors and regulatory authorities has been established, such JSON structures (but not as "file") can be used, using RESTful web services, to populate the regulatory submission, even when the study is still in its starting phase and only a minor amount of the data has been collected.

It is of course a feature that we would LOVE to implement in future ...