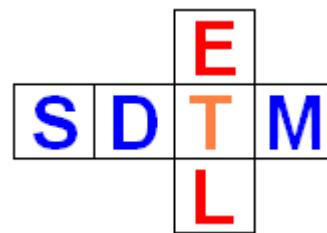# SDTM-ETL 3.1 User Manual and Tutorial

Author: Jozef Aerts, XML4Pharma

Last update: 2014-07-14

## CDISC Dataset-XML Generation

Sofar we have generated SDTM datasets using the option „Generate Transformation (XSLT) Code for SAS-XPT) although we did not create SAS datasets yet, and only generated the SDTM tables with the software itself.
In this part of the tutorial we will demonstrate the different output options which are:
- output of the SDTM datasets in the new CDISC Dataset-XML format
- output of the SDTM datasets in the classic SAS-XPORT (SAS Transport 5) format

The CDISC Dataset-XML format for electronic submissions has been developed to overcome all of the limitations of SAS Transport 5, such as:
- 8-character limitation for variable names and test codes
- 40-character limitation for variable labels
- 200-character limitation for variable values (forcing to split values over records)
- limited amount of publicly (and free) tools available
- not really vendor-neutral[1]

These limitations do not exist in Dataset-XML. Furthermore the new format has a number of additional advantages, such as:
- Perfect alignment with define.xml (both use the same base standard)
- Supplemental variables can remain in the parent domain, they just are flagged as being „non-standard" in the define.xml
- It is not necessary anymore to have a COVAL, COVAL1, etc. variables in case there are comments of over 200 characters (which is often the case)

And there are some very interesting opportunities:
- incorporation of audit records on data points or records
- incorporation of data points from electronic health records (such as using HL7-v3 or HL7-FHIR as format), as Dataset-XML is extensible.

For further information, see the slides presented at the European CDISC Interchange 2014 (Paris).

If you do not have the DM and SV mappings loaded, do now, using „File – load Study define.xml" (or use CTRL-D). After having done so, you can again filter the visible domains (see previous part of the tutorial „Creating mappings for the SV (Subject Visits) domain") using „View – View/Hide Domains". This e.g. results in:

---

1 Although the specification of the format has been published (the TS-140) document, it is extremely difficult implement in modern software that use freely available software languages (Java, C#, …). As such, new vendors start with a serious competition disadvantage.

with the rows „CES:DM" and „CES:SV" containing the mappings for the DM and SV domain respectively, and the row „CES:GLOBAL" containing mappings for global (reusable) variables such as RFSTDTC (Reference Start Date/Time)[2].

Variables that are displayed in grey do already contain a mapping. Variables in red are „required", variables in blue „expected" and variables in green „permissible".

We will now again execute the mappings on our test clinical data. To do so, use the menu „Transform". This results in:
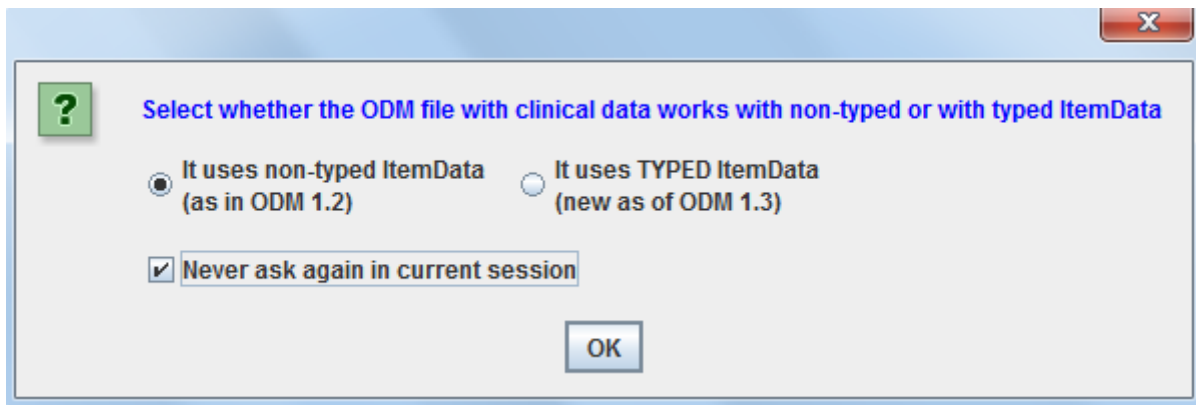


We have the choice between generating the transformation code for generation of SAS Transport 5 (SAS-XPT) datasets, or for the new (modern) CDISC SDS-XML formatted datasets. We choose the latter option (SDS-XML):
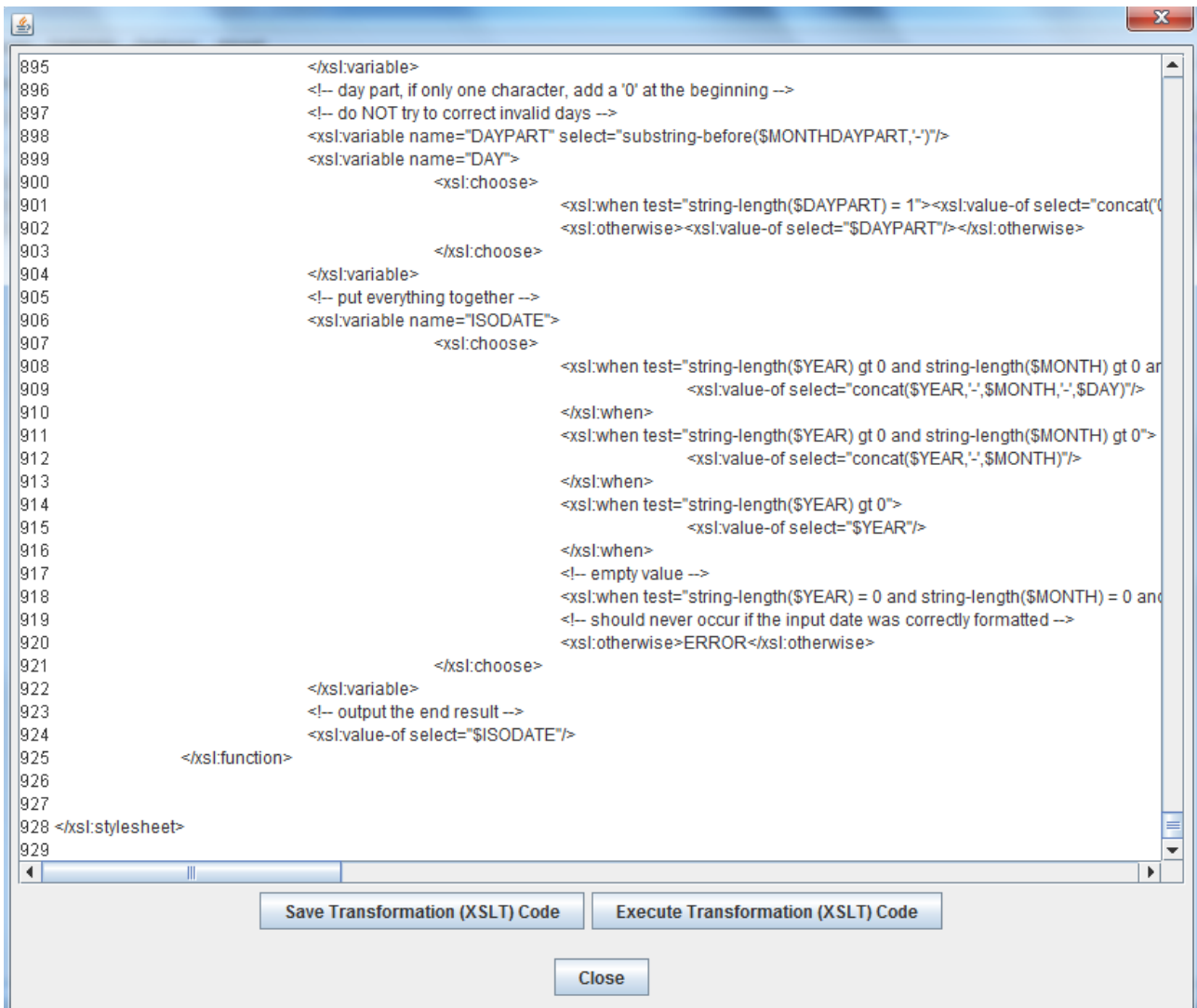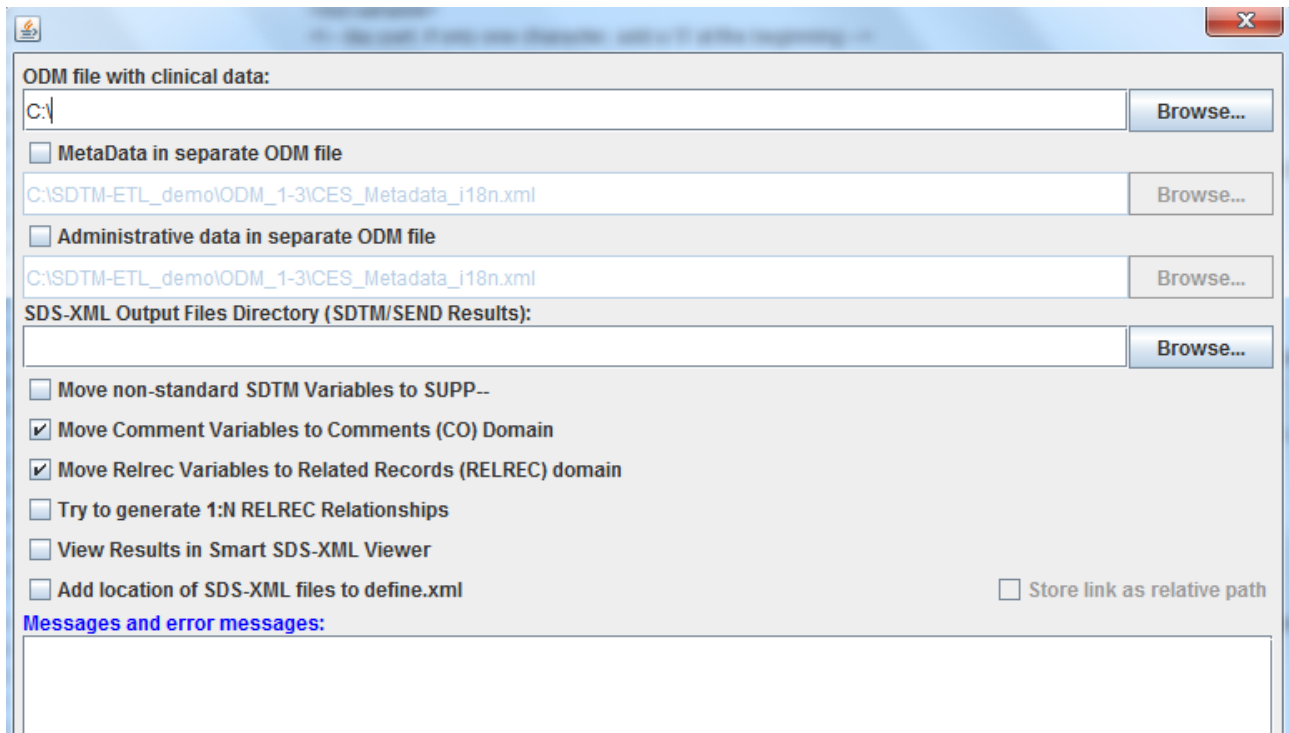


This results in the following dialog:

---

2   Later we will also add RFENDTC (Reference End Date(Time)

We are asked whether our clinical data are „untyped" (classic, using „ItemData" with a „Value" attribute) or „typed" (i.e. ising e.g. ItemDataDateTime). If you are unsure, please ask your EDC vendor or just open an ODM file with clinical data – you will immediately find out.
In most cases, you will avoid that this question pops up again, so check the checkbutton „Never ask again in current session". You are then asked to confirm your choice, and the following dialog is then displayed:



The mapping code is displayed in XSLT (nice for debugging) and the dialog allows you to save the XSLT code for offline execution. You can however also immediately executing the mapping. To do so, click the button „Execute Transformation (XSLT) Code". This results in a new dialog:
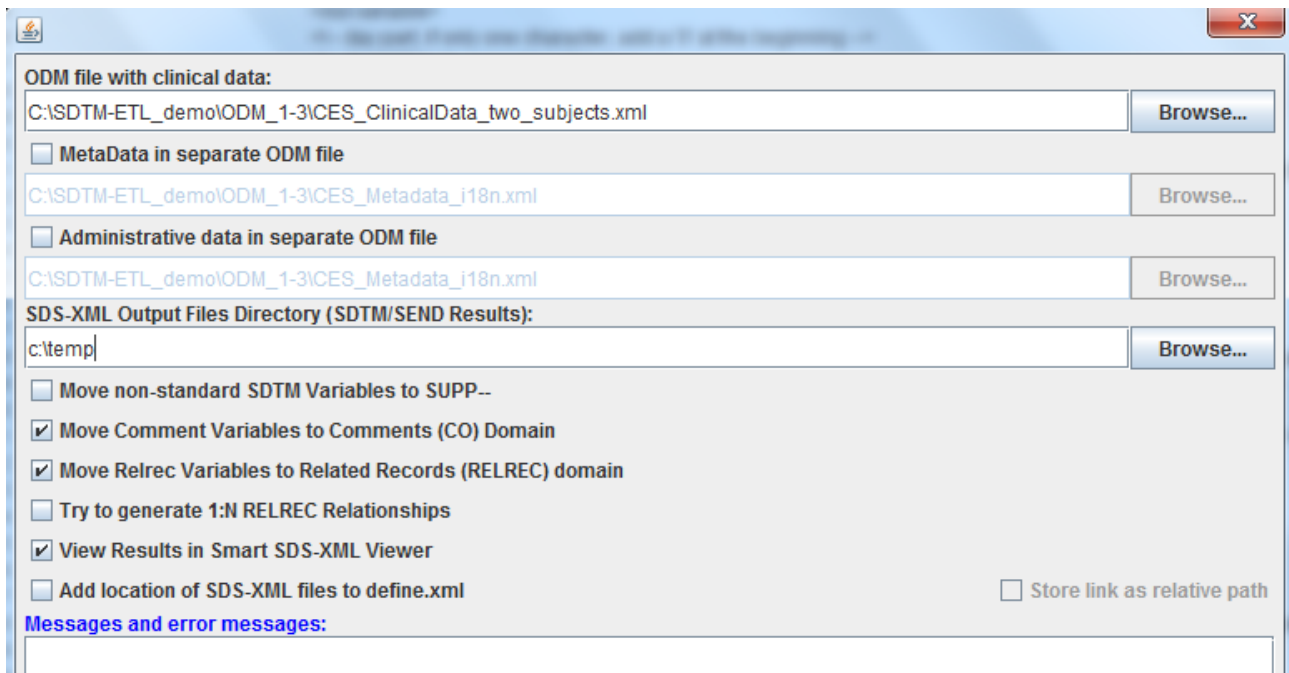
This dialog is very similar to the one obtained when wanting to use SAS-XPT as output, but with a few important differences:

- There is no checkbox anymore asking whether you would want to split values longer than 200 characters – in SDS-XML this is not necessary anymore
- The checkbox „View Result SDTM Tables" has been replaced by a checkbox „View Results in Smart SDS-XML Viewer". We will later see that this viewer is considerable „smarter" than the SASViewer used for inspecting SAS-XPT files

The checkbox „Move non-standard SDTM variables to SUPP--" is still available, but you are discouraged to check it, as keeping the „supplemental" variables in the parent domain is exactly one of the advantages of SDS-XML[3].

At the moment, you are still encouraged to move „Comment" variables into a separate „CO" dataset, although also this might change in future.
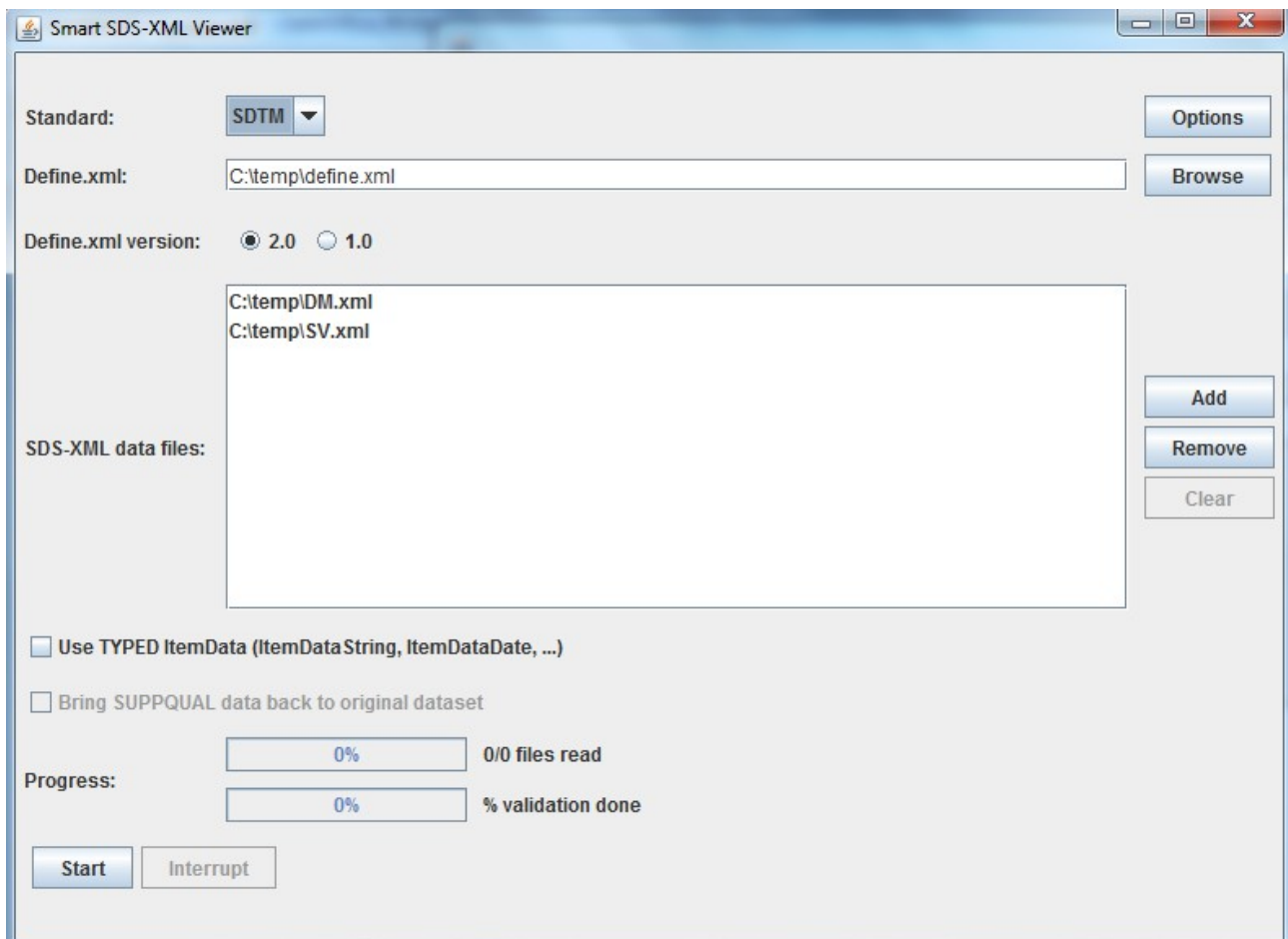
---

3   It still has to be comfirmed by the FDA that they will accept submissions in the new format in the near future.

Now also check the checkbox „View results in Smart SDS-XML Viewer":

I also already added the location of the file with ODM clinical data and of a directory where the SDS-XML files will be written to.
When now clicking the „Execute Transformation on Clinical Data", the transformation executes and after a few seconds the „Smart SDS-XML Viewer" pops up:

A define.xml file has also be generated and written into the same directory as where the generated SDS-XML files reside.

You can now use all possible options of the „Smart SDS-XML Viewer". These are explained in a separate tutorial that you can download from the Open-Source „SourceForge" website:

http://sourceforge.net/projects/smart-sds-xml-viewer/

Also remark that the „Smart SDS-XML Viewer" is open source[4], so you can download the source code, extend it, change it, do anything you want to do with it. You can even replace the version that comes with SDTM-ETL with your own (e.g. extended) one, by replacing the files in the SDTM-ETL directory „Smart_SDS-XML_Viewer_software" (of course everything on your own responsibility!).

Starting the „Smart SDS-XML Viewer" in our case results in:
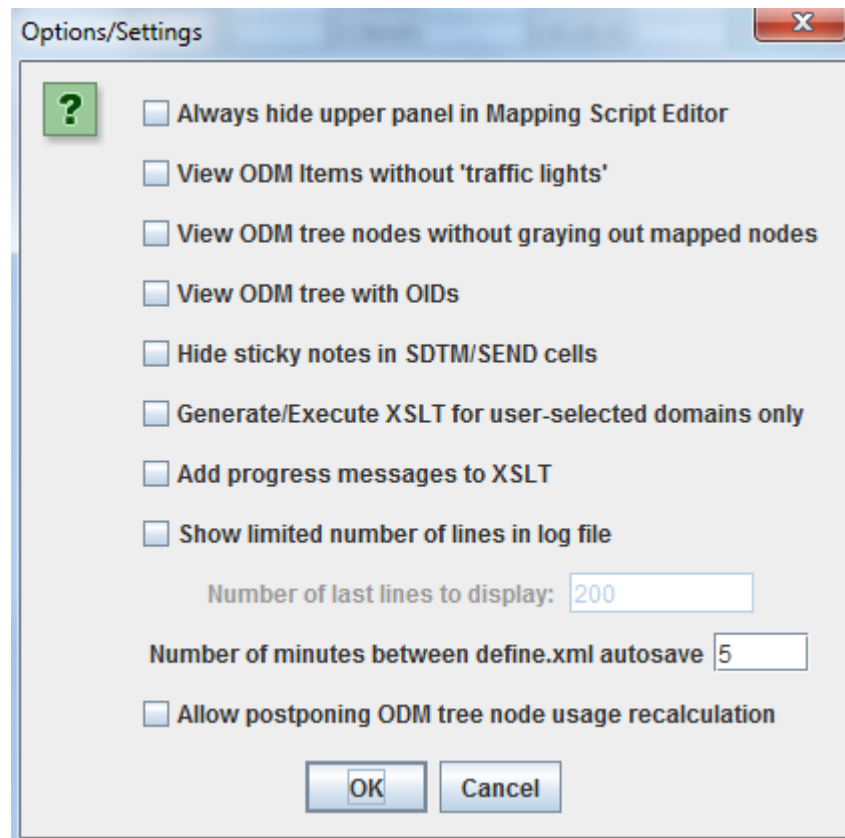


and for the SV dataset:



Remark that unlike with the SASViewer, the „Smart SDS-XML Viewer" also displays the metadata (like the define.xml datatype) when hovering the mouse over the column header.

Much more is however possible with the „Smart SDS-XML Viewer" and you are invited to download the tutorial from the SourceForge site, but you can also look at a YouTuve demo video at:
https://www.youtube.com/watch?v=wRBktVtB47s&noredirect=1

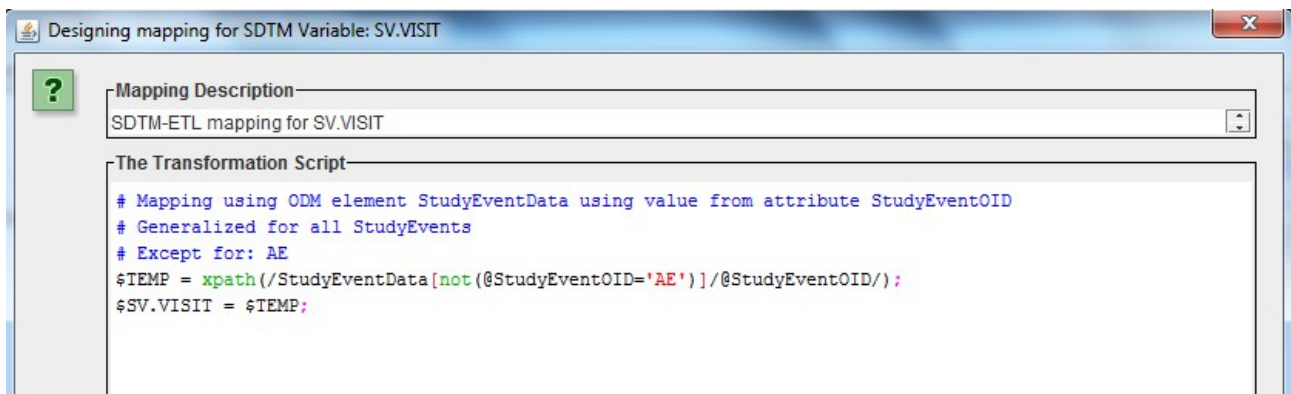Some options and settings for SDTM-ETL

There is a good number of options you can choose from when working with SDTM-ETL. You can change them using the menu „Options – Settings":

---

4  Once the CDISC specification is final, the source code will be made available through SourceForge.
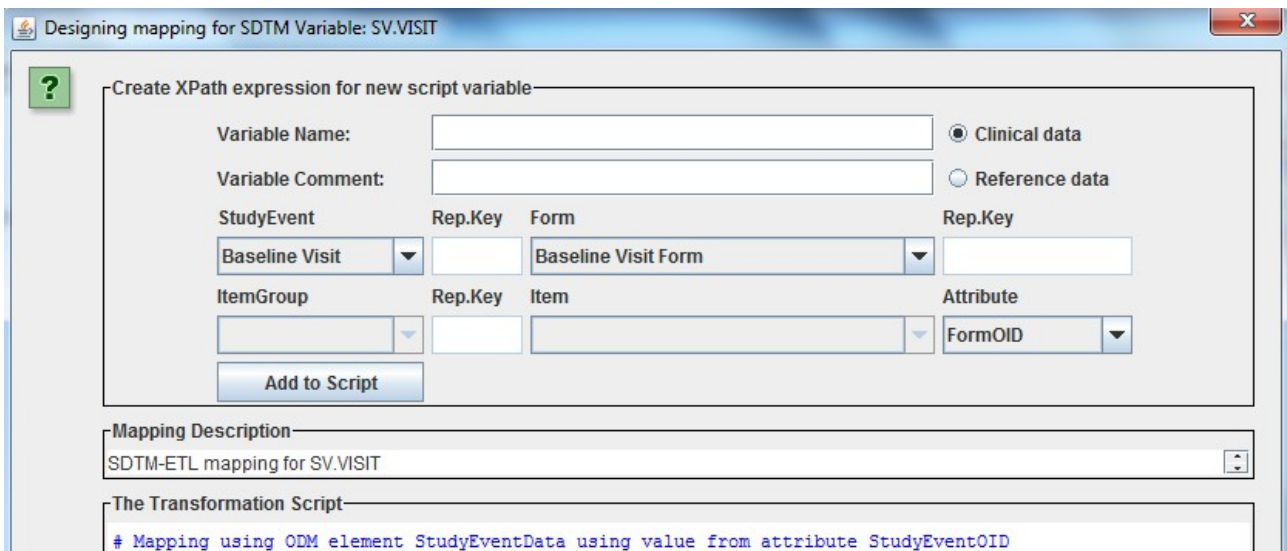
The number of options is steadily increased, so you might see a slightly different picture with your distribution.
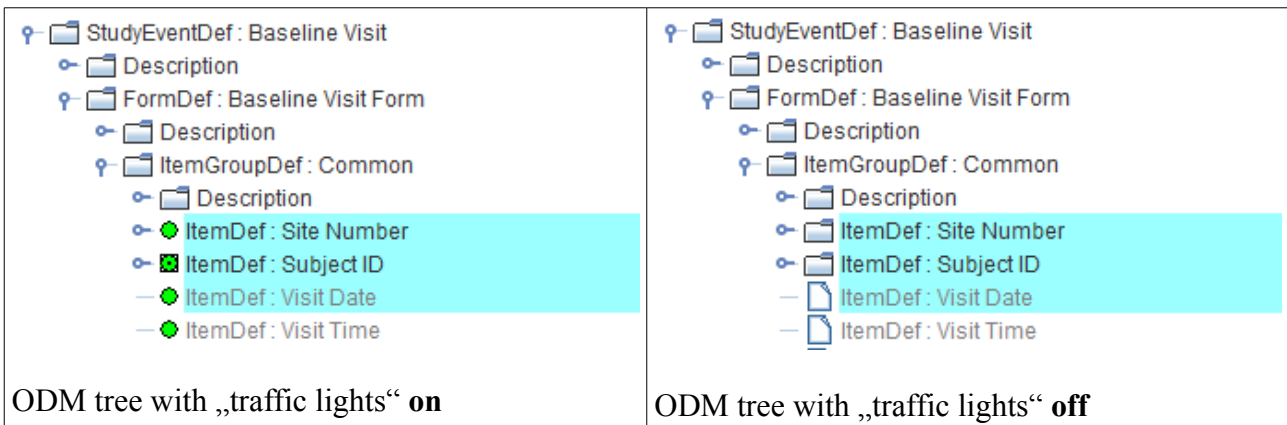
The first option „Always hide upper panel in Mapping Script Editor" always hides the upper panel in the mapping editor. So when it is checked, and editing a mapping, you will see:



instead of:

The second option „View ODM Items without „traffic lights" allows you to switch off the „traffic lights" in the ODM tree.



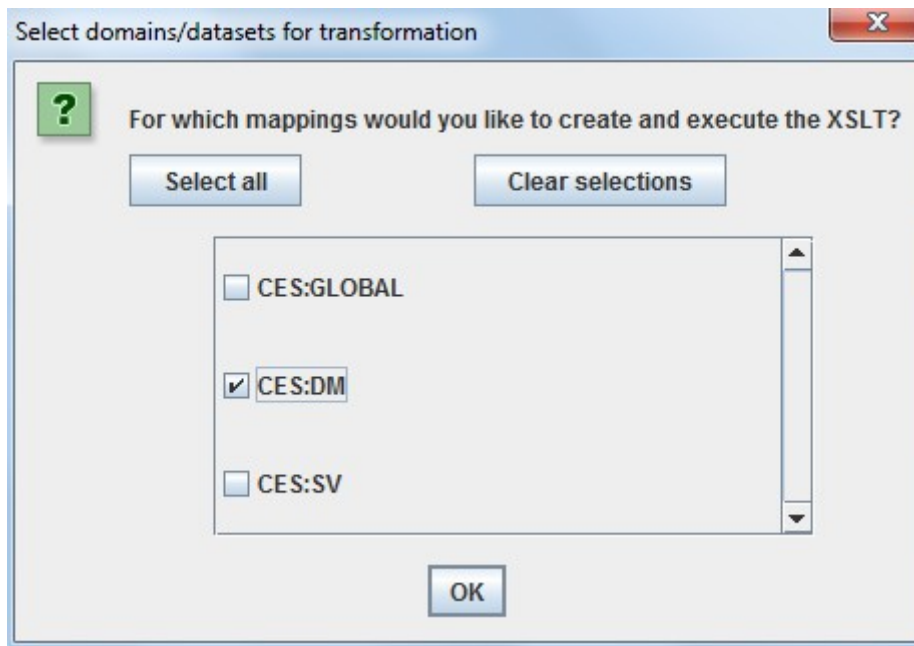| ODM tree with „traffic lights" **on** | ODM tree with „traffic lights" **off** |

Some people love the „traffic lights", other hate them.
Anyway, it is advised to switch them off, in case you see that your system slows down considerably when selecting another SDTM cell. This can happen when you have a very large or complicated ODM tree, as each time that a new SDTM cell is selected, the „traffic lights" need to be recalculated.

The third option „View ODM tree nodes without graying out mapped nodes" is similar.
You might have noticed that when an Item in the ODM tree has been used in a mapping, it is „grayed out". Also here, some users like this feature, others don't.

The next option „Hide sticky notes" will be explained later, when we discuss the feature of adding „sticky notes" to SDTM variables.

A new option as of SDTM-ETL 3.0 is the option „Generate/Execute XSLT for selected domains only", and is very useful, especially when you have mappings for several domains, but only want to test the execution of the mappings for a single or a few selected domains only. In case this checkbox is checked, each time you want to execute a mapping, a dialog will be displayed asking you for which of the current domains you would like to execute the mappings. For example you may then see:

where it is requested that only the mappings for the DM domain are executed.
Selecting a single of many domains considerably increases the speed of execution, so this options can be very useful.

We now continue with the generation of the XSLT for generating Dataset-XML submission files.
The XSLT is generated and presented to the user:

```
1929                          </xsl:variable>
1930                          <!-- day part, if only one character, add a '0' at the beginning -->
1931                          <!-- do NOT try to correct invalid days -->
1932                          <xsl:variable name="DAYPART" select="substring-before($MONTHDAYPART,'-')"/>
1933                          <xsl:variable name="DAY">
1934                                  <xsl:choose>
1935                                          <xsl:when test="string-length($DAYPART) = 1"><xsl:value-of select="concat('0
1936                                          <xsl:otherwise><xsl:value-of select="$DAYPART"/></xsl:otherwise>
1937                                  </xsl:choose>
1938                          </xsl:variable>
1939                          <!-- put everything together -->
1940                          <xsl:variable name="ISODATE">
1941                                  <xsl:choose>
1942                                          <xsl:when test="string-length($YEAR) gt 0 and string-length($MONTH) gt 0 ar
1943                                                  <xsl:value-of select="concat($YEAR,'-',$MONTH,'-',$DAY)"/>
1944                                          </xsl:when>
1945                                          <xsl:when test="string-length($YEAR) gt 0 and string-length($MONTH) gt 0">
1946                                                  <xsl:value-of select="concat($YEAR,'-',$MONTH)"/>
1947                                          </xsl:when>
1948                                          <xsl:when test="string-length($YEAR) gt 0">
1949                                                  <xsl:value-of select="$YEAR"/>
1950                                          </xsl:when>
1951                                          <!-- empty value -->
1952                                          <xsl:when test="string-length($YEAR) = 0 and string-length($MONTH) = 0 and
1953                                          <!-- should never occur if the input date was correctly formatted -->
1954                                          <xsl:otherwise>ERROR</xsl:otherwise>
1955                                  </xsl:choose>
1956                          </xsl:variable>
1957                          <!-- output the end result -->
1958                          <xsl:value-of select="$ISODATE"/>
1959                  </xsl:function>
1960
1961
1962 </xsl:stylesheet>
1963
```
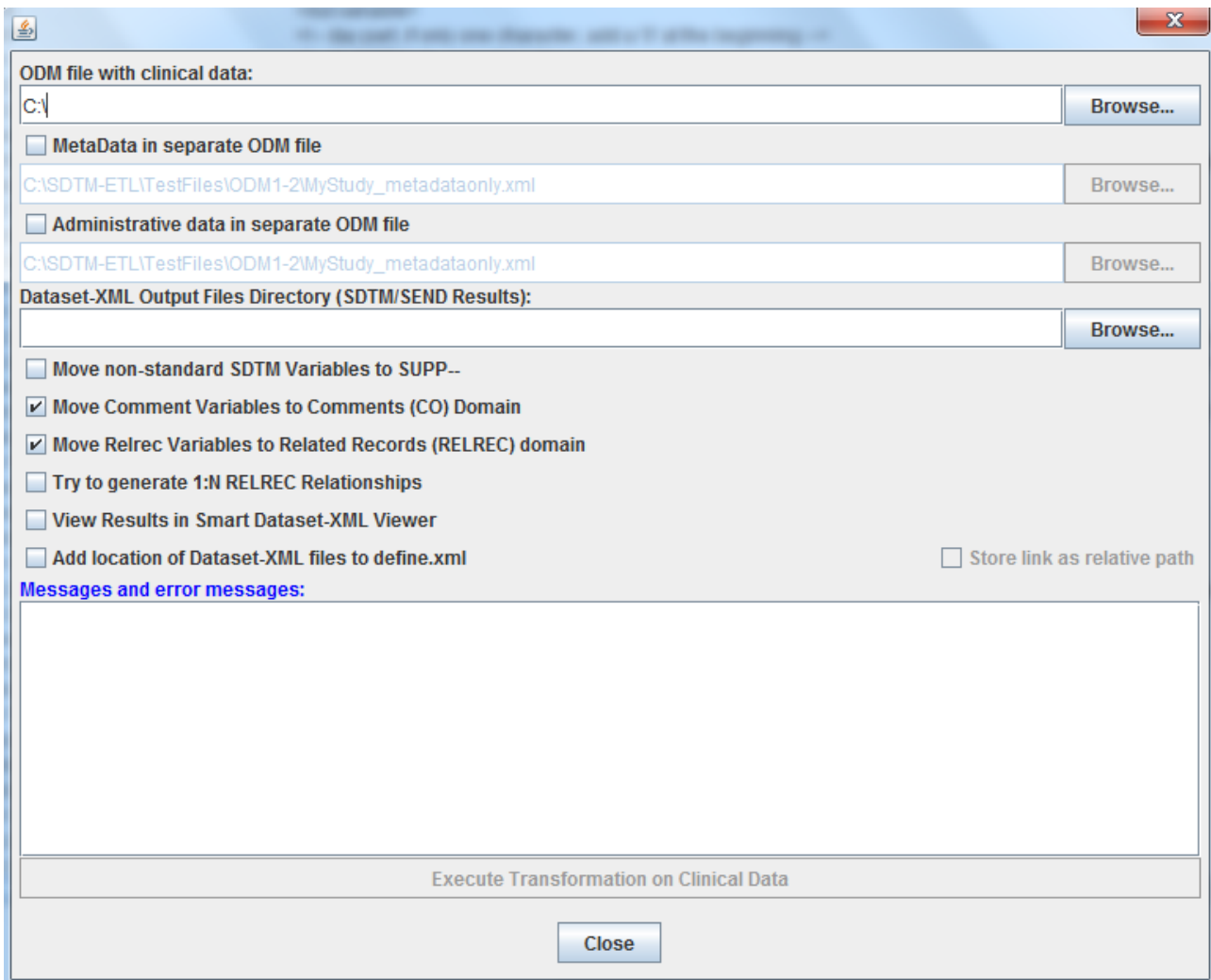
Save Transformation (XSLT) Code        Execute Transformation (XSLT) Code
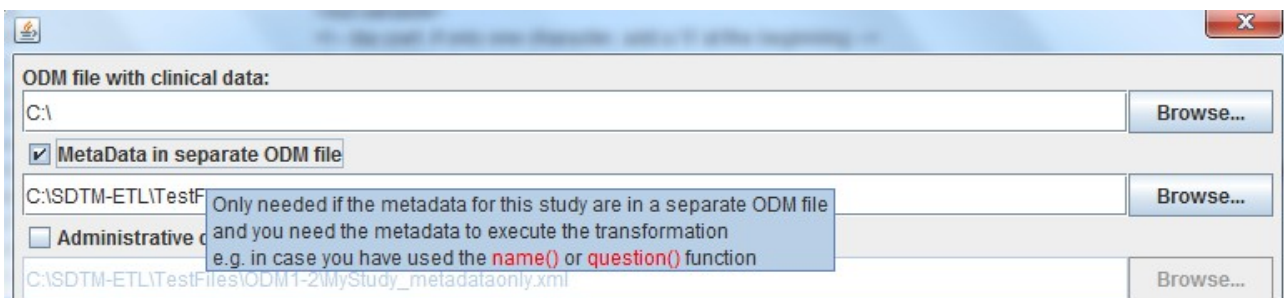
Close

No worry! You do not need to understand XSLT to execute the mappings! However, some of our expert users want to see it and some want to save it to disc for offline execution.

Now just click the „Execute Transformation (XSLT) Code" button. The following dialog is now displayed:

You are requested to locate the file with the clinical data.

If you use special functions in your script that requires the metadata (like the „name()" or „question()" function, and your file with clinical data does <u>not</u> contain the metadata, you also need to select the checkbox „", and provide the location of the file with the metadata (study design):



Also, if you need to retrieve some information from your „AdminData" (e.g. site information), and this data is in a separate file, you also need to check the checkbox „" and provide the location of the file with administrative data:

Now also provide the location of an existing directory where the output of the transformation (the Dataset-XML files) will be written to, e.g.:



You can now also decide whether non-standard variables should be moved to SUPP-- datasets (which is unfortunately still required by the FDA). Also, if you want to move „comment" variables to a CO dataset (also still required by SDTM), check the corresponding checkbox. Similarly for RELREC records. You can also select whether the system should try to generate 1:N RELREC records. This will be treated in detail later.

Finally, if you want to see the results of your transformation using the open source „Smart Dataset-XML Viewer"[5], check the corresponding checkbox.

The checkbox „Add location of Dataset-XML files to define.xml" will take care of writing the def:leaf elements in the define.xml file. This is only required in the final stages of developing a submission.

Based on the current (July 2014) requirements of the SDTM and the FDA, the typical set iof options will look like:



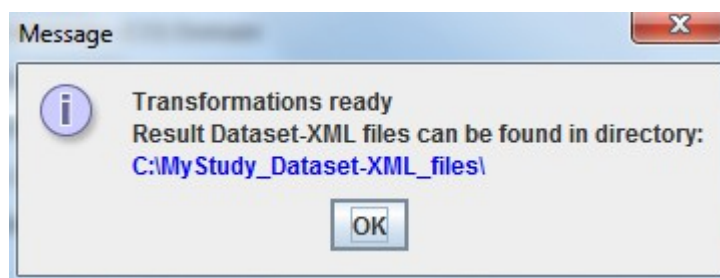Now it is time to click the „Execute Transformation on Clinical Data" button:

---

5   The „Smart Dataset-XML Viewer" comes as a part of the „SDTM-ETL" software. The system has however been designed in such a way, that you can replace an older version of the viewer by a more recent one, or even by your own one (as the „Smart Datase-XML Viewer" is open source, you might have extended it).
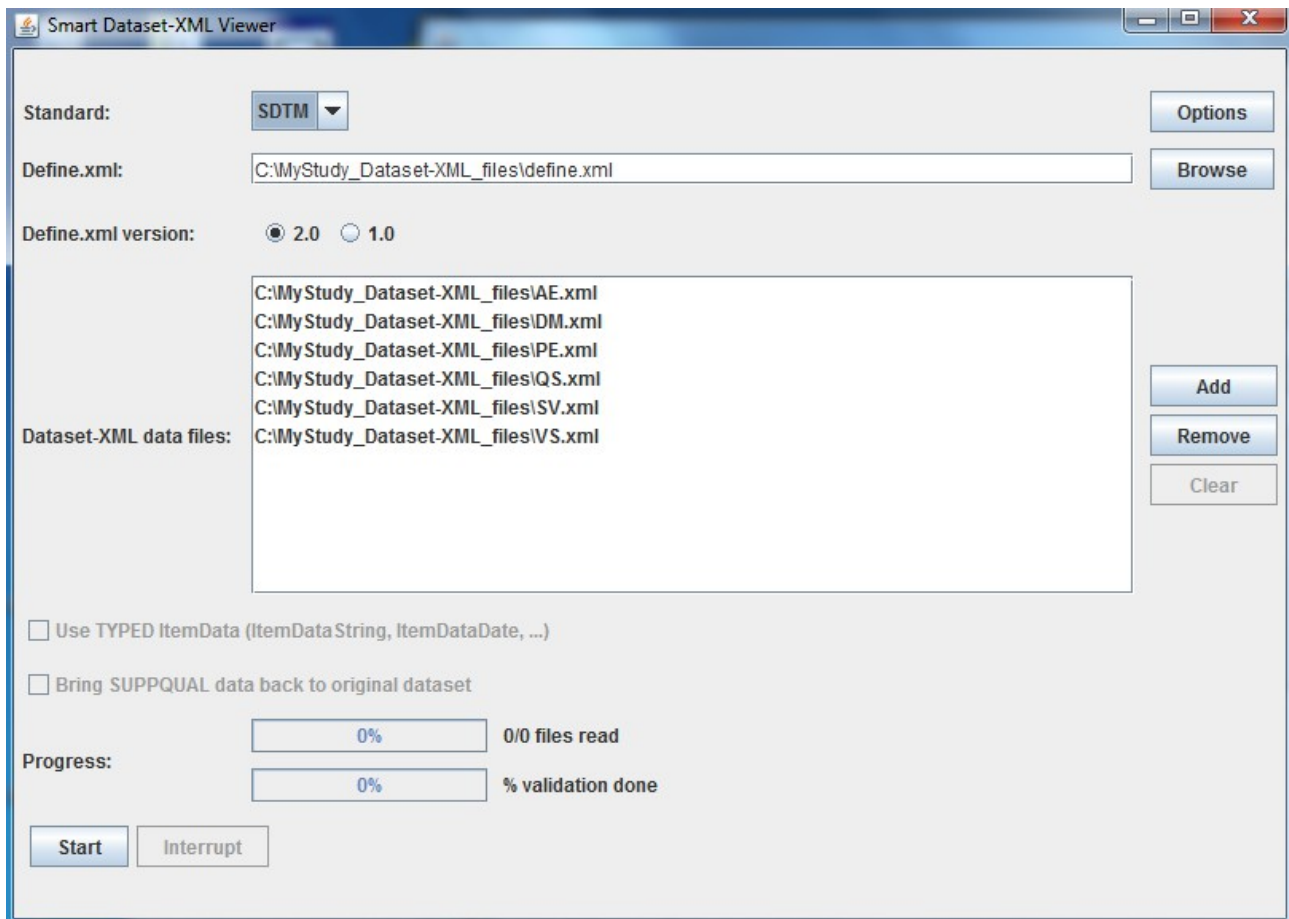
Some messages might show up in the „messages and error messages" console. For example, if you have used the option in „Options - „ also a message will be printed each time a new dataset has been generated. This console will also show you all the messages from your „print" statements in your mappings, which is very useful when debugging mappings.
For example:



When everything is finished, the following message is displayed:



and if you choose to use the „Smart Dataset-XML Viewer" for inspecting the results, a new window is displayed already containing all the necessary information, e.g.:

and after clicking the „Start" button resulting in:



For more information about working with the „Smart Dataset-XML Viewer", please download the manual / tutorial from the Sourceforge site http://sourceforge.net/projects/smart-sds-xml-viewer/files/.
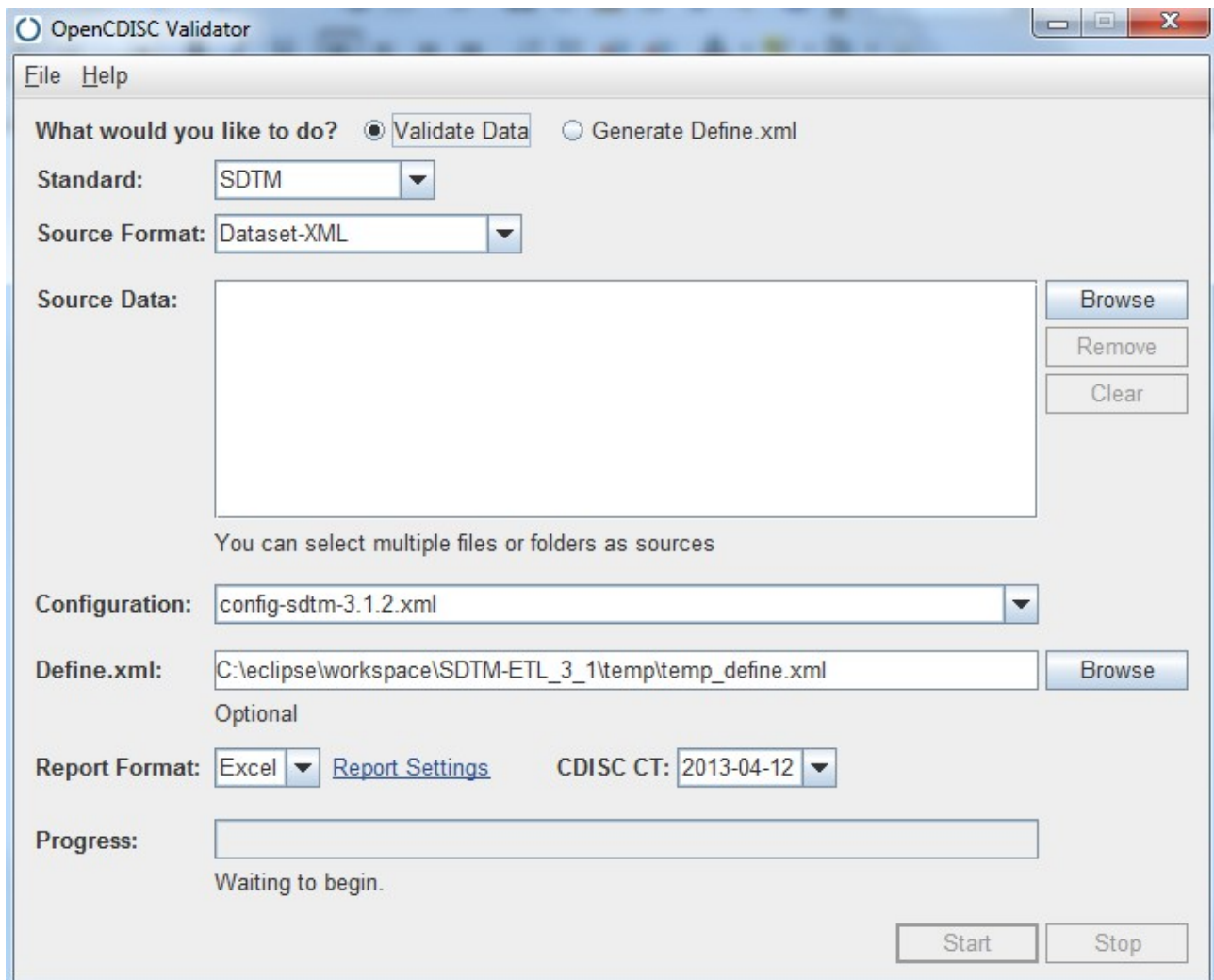
After inspection, you can close all but the main window, and continue with the development of new mappings.

Validating your generated Dataset-XML files using OpenCDISC

The newest version of OpenCDISC (v.1.5) supports Dataset-XML and has been integrated into the SDTM-ETL software[6]. So if you use the menu „Validate – Validate SDTM Dataset-XML Records against define.xml":
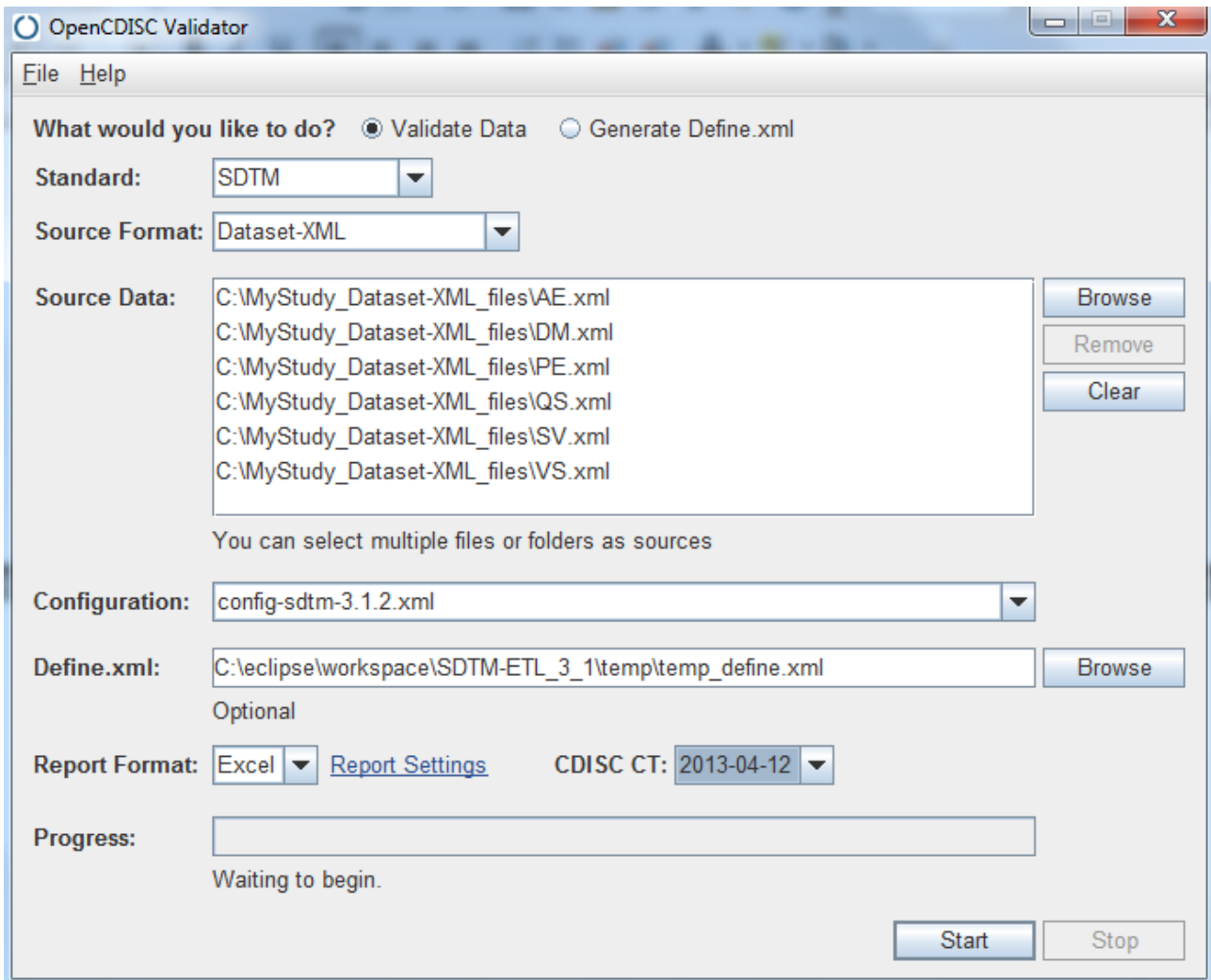


OpenCDISC will be started automatically and most of the necessary fields will already be filled:
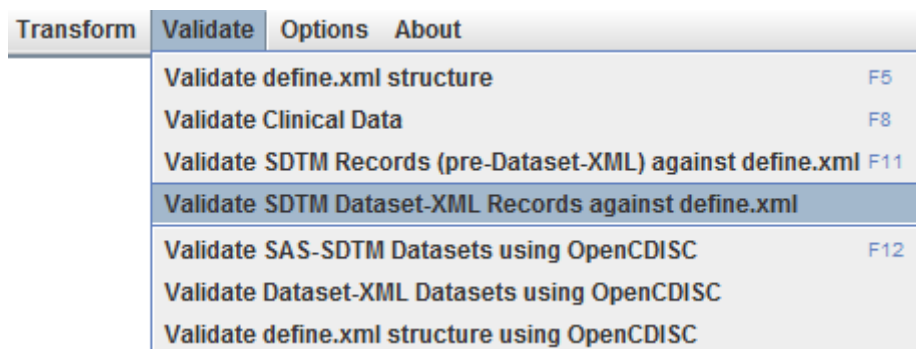


You will still need to provide a list of generated Dataset-XML files, e.g. the ones we created in the previous step:

6   Remark that we found some bugs in the Dataset-XML implementation which we have corrected in the SDTM-ETL integration. These bugs have been reported to OpenCDISC but no formal update has been provided by OpenCDISC since.
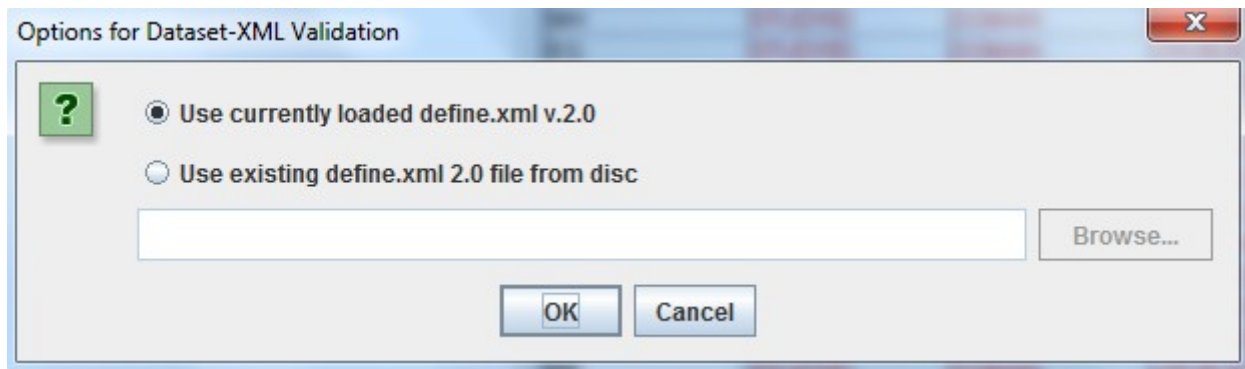
For more information of validating Dataset-XML files using OpenCDISC, please refer to the
OpenCDISC website. Please also note that the validation rules of OpenCDISC is the „OpenCDISC
interpretation" of the SDTM/SEND standard, so not the XML4Pharma interpretation and even not
the „CDISC interpretation". So if you get validation errors that you do not understand, or disagree
with, please use the „OpenCDISC Forum" to discuss these with the OpenCDISC people.
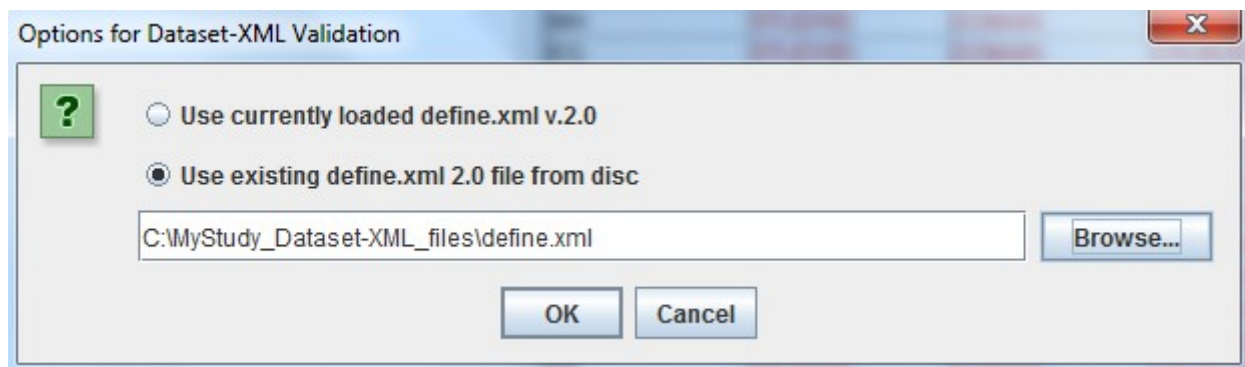
For validating the generated Dataset-XML files, you can also use the menu „Validate -
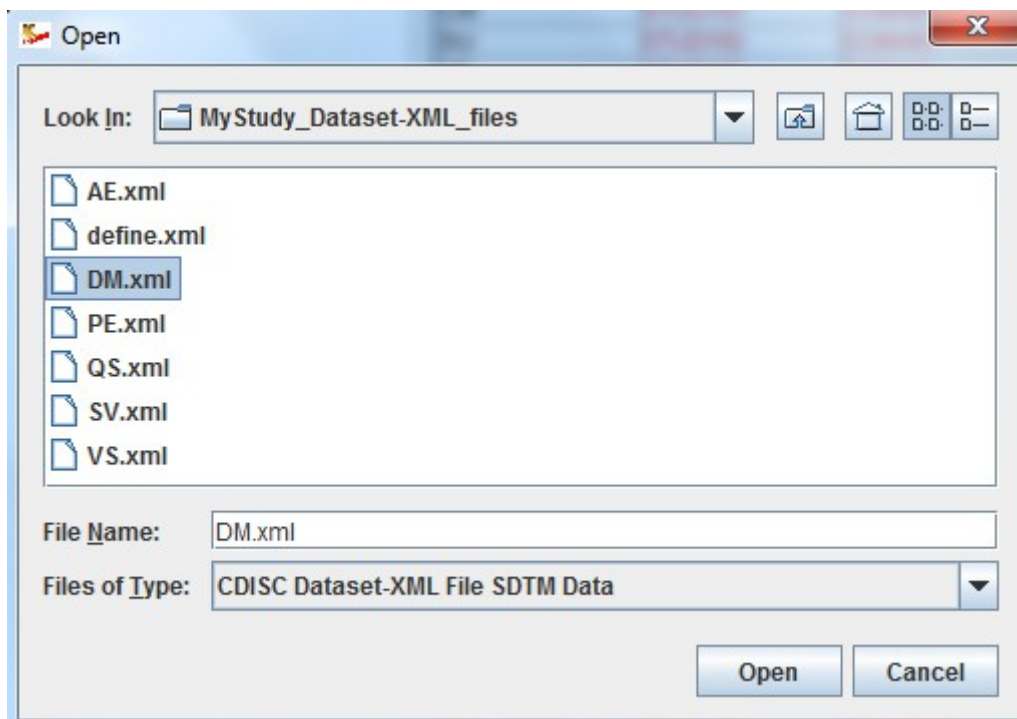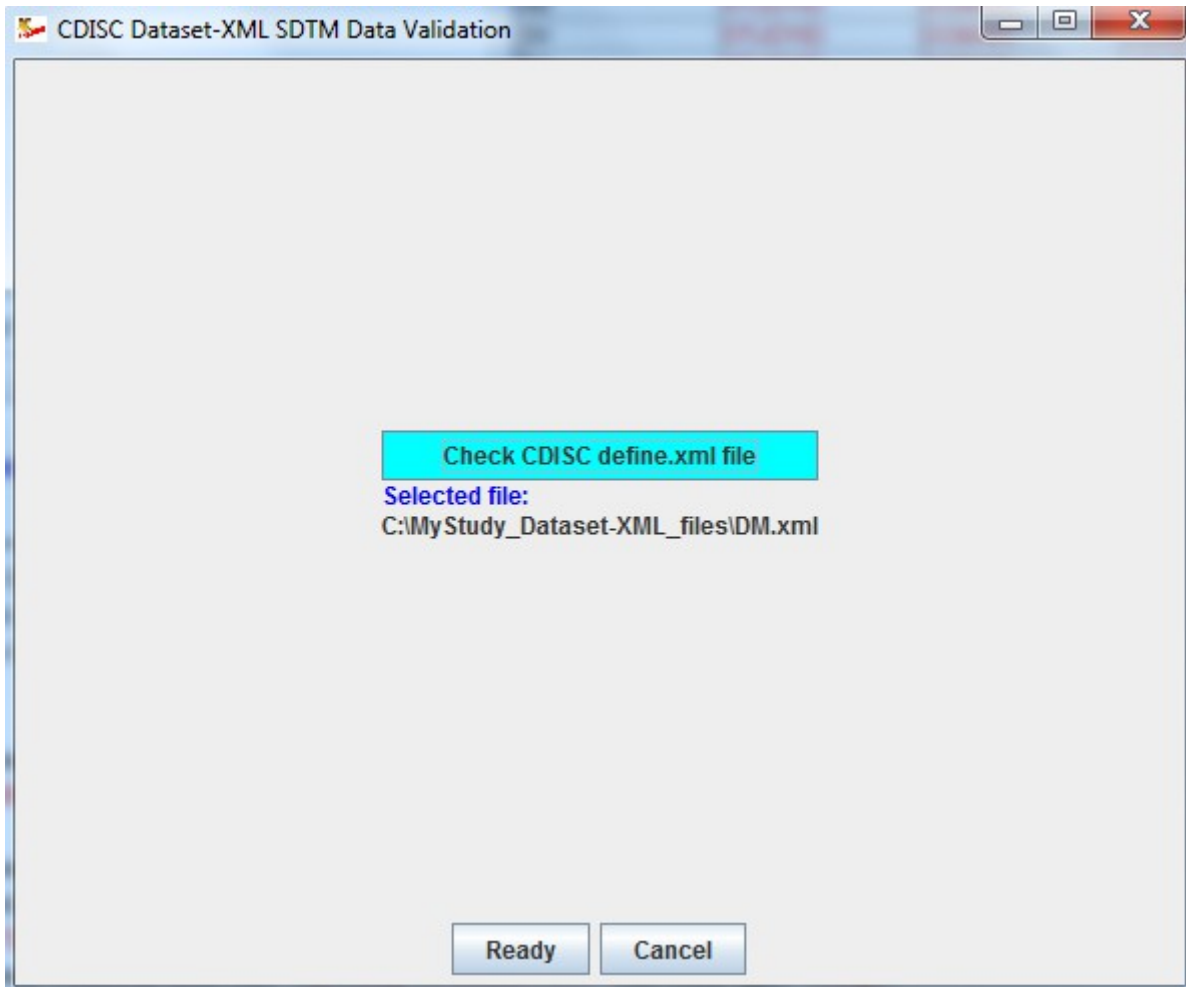


The following dialog is displayed:

You can either use the currently loaded define.xml (a copy will be made by the software), or if you do already have one (as in our case), you can select to use the one that you have already on disc:



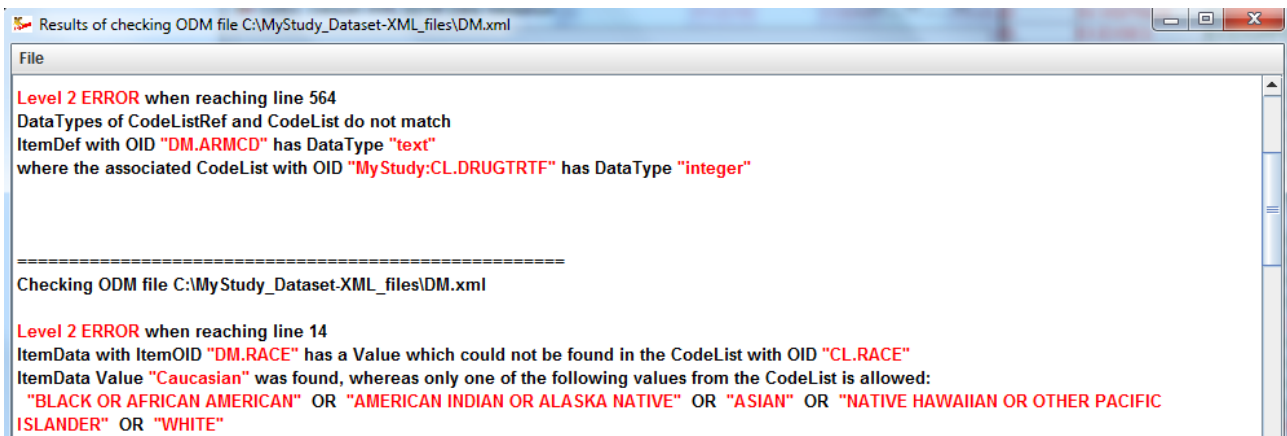after clicking OK, you can start selecting a Dataset-XML file, e.g.:



Unlike with OpenCDISC, you can only validate one Dataset-XML file at the time, leading to the following dialog:

TODO: text on button is not correct.

And e.g. as a result of the validation:



There is an important difference in validation approach between the OpenCDISC validator and the one from SDTM-ETL:
Essentially, OpenCDISC validates the contents of the files against the SDTM rules and official CDISC controlled terminology. For example, if you developed your own codelist for an SDTM variable (which is stored in the define.xml) instead of using the CDISC one, OpenCDISC will report an error.
The SDTM-ETL validator for Dataset-XML files however only validates the data against the metadata in the define.xml. So it does not try to interpret the SDTM-IG, nor does it do a look up in

any published CDISC-CT. Leading in the validation is only information in the define.xml.
This does thus also mean that the SDTM-ETL validator for Dataset-XML does not e.g. validate that
an end date/time (--ENDTC) does not come before the corresponding start date/time (--STDTC).

So validating your Dataset-XML files with OpenCDISC and with the build-in validator WILL give
different results, which should be seen as complimentary.

Generating SAS-XPT (Transport 5) SDTM files

You can still generate SAS Transport 5 SDTM files using SDTM-ETL 3.0/3.1. This is probably the safe option as long as there is no official announcement that the FDA accepts submissions in the new Dataset-XML format (or it must be that you take part in the FDA pilot project).

The process for generating SDTM/SEND datasets is very similar to what is described here (you will see many of the dialogs already displayed here), and is explained in a separate document.