Smart Submission Dataset Viewer User Manual

By XML4Pharma Last update: 2024-12-16

This user manual and also other manuals and tutorials can always be found on our website: http://www.xml4pharma.com/Smart_Submission_Dataset_Viewer/

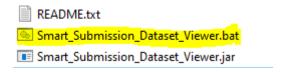
Table of contents

Table of contents	l
Starting the application	1
Selecting data standard, define.xml and to-be-loaded files	2
Viewing the results - basic features	
Viewing values for Supplemental Qualifiers	11
Looking up basic subject information in DM or ADSL	16
Sorting and filtering	17
Moving columns around and hiding columns	24
Connecting Related records (RELREC)	26
Exporting the table	30
Basic validation	31
Search and View	32
Filtering before loading Setting Viewer Options Extended validation	36
Setting Viewer Options	40
Extended validation	47
Validation using CDISC CORE	48
Using APIs	49
The "properties.dat" file	52

Starting the application

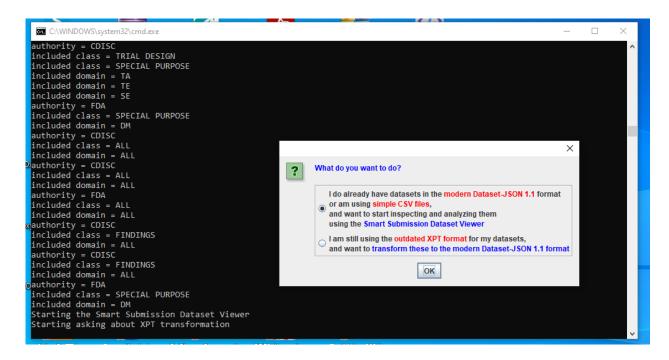
The software comes as a Java application. For installation instructions, please find them <u>here</u>. The document also contains information about setting up a desktop shortcut for the software.

Essentially, one starts the software by a double-click on "Smart_Submission_Dataset_Viewer.bat" (Windows) or "Smart_Submission_Dataset_Viewer.sh" (Linux and Mac) in the folder where you installed the "Smart Submission Dataset Viewer". For example:



or on the corresponding shortcut item on the desktop.

This will start with the opening screen:



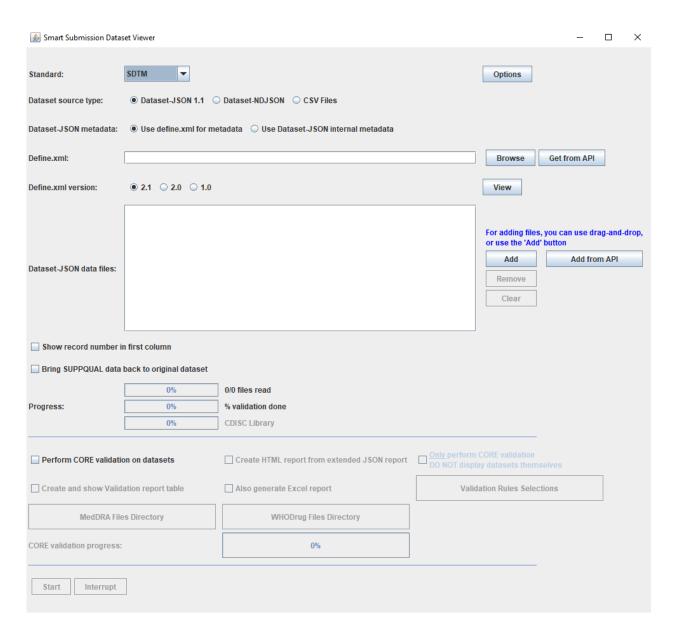
At the same time, a "console window" (the "black" window) is opened, that will display all information messages produced by the software, and which are also stored in a log file, which can be find in the folder "logs".

The opening screen asks whether one already has modern CDISC Dataset-JSON 1.1 files available, or that one wants to convert existing SAS-XPT files to Dataset-JSON. For the latter, see the separate tutorial "Generating Dataset-JSON from SAS-XPT".

If one already uses Dataset-JSON, and one does not want to be asked again about conversion from SAS-XPT, this opening screen can be skipped by setting a parameter in the "properties.dat" file. See the section "The 'properties.dat' file".

Selecting data standard, define.xml and to-be-loaded files

In case no conversion from SAS-XPT file is needed as one already has files in Dataset-JSON format, clicking "OK" leads to the main screen:



First, one will want to select the standard for which the viewer will be used. Current choices are: "SDTM", "SEND", "ADaM" and "Other tabular". In this manual, we will use an SDTM example. These choices are there in order to optimize display and additional features. For example, CDISC ADaM does not have the concept of "Supplemental Qualifier datasets" (SUPPQUAL), so when ADaM is selected, the checkbox "Bring SUPPQUAL data back to original dataset" will be disabled.

Also, when one selects "ADaM, some new choices appear, such as:



We will explain the "Options" button on the right top later.

After having selected the standard, one must select in which format the data will delivered to the viewer. The choices are:

- Dataset-JSON 1.1: This is the "regular" CDISC Dataset-JSON 1.1 format, for which the specification can be found at: https://cdisc-org.github.io/DataExchange-DatasetJson/doc/dataset-json1-1.html

You will probably want to have your submission data files in this format in 90% of the cases.

Information about software packages and tools to develop the mappings for, and generate submission files in Dataset-JSON can be found at:

https://wiki.cdisc.org/display/DSJSONHACK/Hackathon+I+Projects

Some SDTM/SEND mapping tools, such as the "<u>SDTM-ETL</u>" software, also allow to generate submission datasets in Dataset-JSON format.

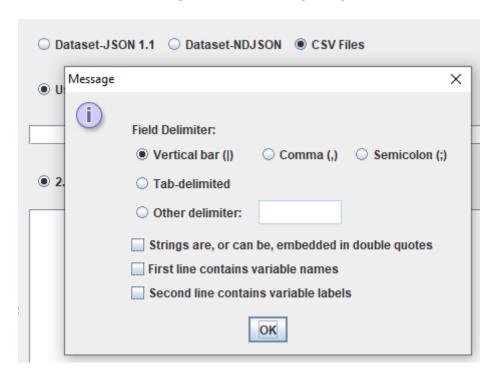
- Dataset-NDJSON

This is a special form of Dataset-JSON 1.1 in which each record is on a separate line in the file, each line containing a separate JSON object (ND stands for "newline delimited"). It can be advantageous to be able to process for very large datasets in software tools that store the JSON in memory.

The specification can also be found on the CDISC Wiki.

The Smart Submission Dataset Viewer however uses "streaming" when reading the contents of the files, so that, at least for our software, there is no difference in performance between Dataset-JSON and Dataset-NDJSON. But as people might want to generate Dataset-NDJSON for other reasons, we of course need to support it.

- CSV (Comma-separated Values). When selecting "CSV", the following dialog is shown:



allowing the user to select or assign the delimiter between fields to be used and further details how the data is organized.

Using CSV can be useful when e.g. obtaining SDTM/SEND/ADaM files as Excel or any other spreadsheet format.

Suppose we selected "SDTM" and "Dataset-JSON", then we need to select whether we want to use a define.xml file for loading the metadata for the files, or that we just want to use the metadata that is already present in the Dataset-JSON files themselves.

Dataset source type:	Dataset-JSON 1.1	et-NDJSON CSV Files
Dataset-JSON metadata:	Use define.xml for metadata	O Use Dataset-J SON internal metadata

As the amount of metadata in the Dataset-JSON files themselves is very limited, it is always best to use the define.xml for loading the metadata. This however requires to have a good quality define.xml. Generating such a define.xml starting from an Excel file (as a "specification") is not always a good idea, as it often produces low-quality define.xml files. There are however different mapping tools on the market that generate a quality define.xml "synchronously" with the mappings themselves. There is also a good number of software tools for developing, generating and updating/extending define.xml on the market that come with a user-friendly graphical user interface.

Depending on whether one selects "Use define.xml for metadata" or "Use Dataset-JSON internal metadata", some features will become available or disabled. For example, the feature "Bring SUPPQUAL data back to original dataset" requires information from the define.xml. As such, this checkbox will be disabled when one selected "Use Dataset-JSON internal metadata". Also many data validation features, such as checking values against their variable codelist, are only possible when using the define.xml for loading the metadata (see later).

When having selected "Use define.xml for metadata", one should now first select the define.xml, by using the "Browse" button, and then select the define.xml that is applicable to the datasets to be loaded.

Dataset-JSON metadata:	Use define.xml for metadata	○ Use Dataset-J SON internal metadata	
Define.xml:			Browse Get from API
Define.xml version:	● 2.1 ○ 2.0 ○ 1.0		View

In the current version of the software, one also sees a "Get from API" button. This is currently for demonstration purposes only, and is explained in the section "Using APIs".

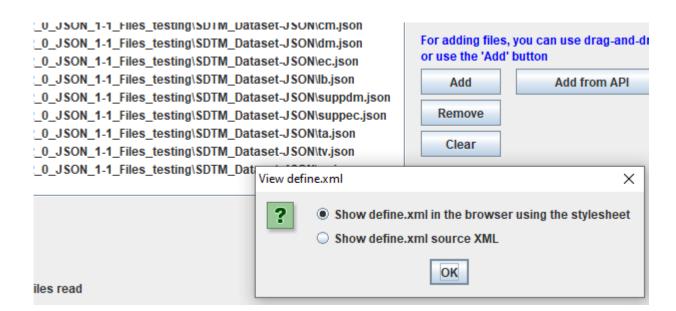
In this manual, we will explain many of the features using the Dataset-JSON 1.1 files from the "SDTM Metadata Submission Guidelines v2.0" published by CDISC.

So we select its define.xml e.g.:

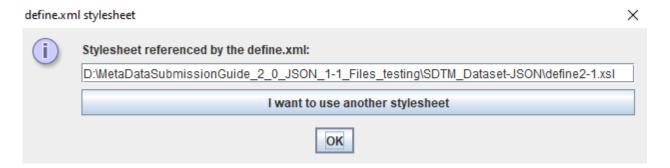


Normally, the software will recognize which version of Define-XML the define.xml file implements, and in this case automatically select "2.1" for "Define.xml" version, but it is always a good idea to check this.

One can now also inspect the contents of the define.xml by clicking the "View" button. The following dialog is then displayed:

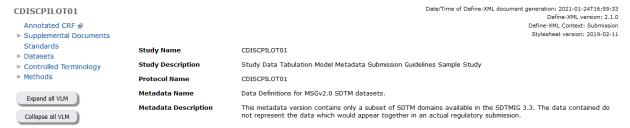


If one chooses "Show define.xml in the browser ...", the system proposes to use a stylesheet it found in the same folder as the define.xml itself, but also allows the user to select another stylesheet:



Remark that the "Smart Submission Dataset Viewer" does not come with Define-XML stylesheets, as the stylesheet is the responsibility of the user's organization.

The stylesheet will then be applied and then display the define.xml as HTML in your default browser, e.g.



Standard's Conformance Notes: 1) The SDTM v1.7/SDTMIG v3.3 datasets were evaluated manually and programmatically by the CDISC SDS MSG Team. At the completion of the SDTM-MSG v2.0, the CDISC SDTM v1.7/SDTMIG v3.3 conformance rules were recently published, but not available by any validation tools to validate. 2) The Define-XML document was evaluated manually and programmatically by the CDISC SDS MSG Team. At th completion of the SDTM-MSG v2.0, the CDISC Define-XML v2.1 conformance rules were not published, nor available by any validation tools to validate. Please ensure that any official regulatory submission of an Define-XML v2.1 document and accompanying data is done in accordance to the respective regulatory health authorities requirements/guidance.

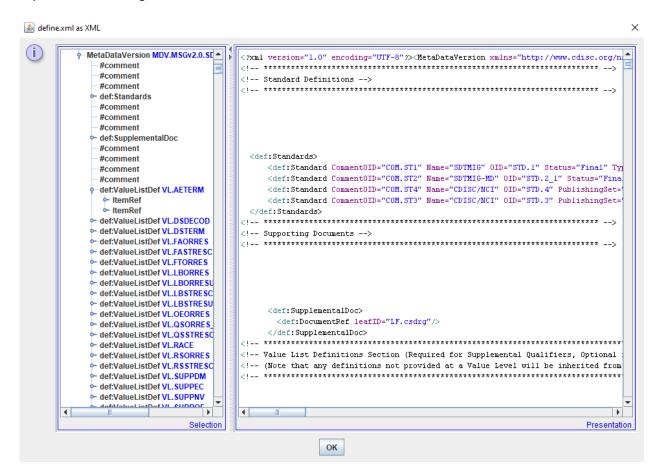
Standards for Study CDISCPILOT01

Standard	Туре	Status	Documentation
SDTMIG 3.3	IG	Final	
SDTMIG-MD 1.1	IG	Final	
CDISC/NCI DEFINE-XML 2020-12-18	СТ	Final	
CDISC/NCI SDTM 2020-12-18	СТ	Final	

Datasets

Dataset	Description	Class	Structure	Purpose	Keys	Documentation	Location
TA [SDTMIG 3.3]	Trial Arms	TRIAL DESIGN	One record per planned Element per Arm	Tabulation	STUDYID, ARMCD, TAETORD		ta.json &
TE [SDTMIG 3.3]	Trial Elements	TRIAL DESIGN	One record per	Tabulation	STUDYID,		te.json @

The second choice "show define.xml source XML" allows to display the define.xml as XML itself in a separate window, e.g.:



where one can quickly navigate to specific sections by clicking on one of the tree nodes in the left part of the window.

This feature is of course only of interest to people who understand the XML structure of Define-XML.

We will then want to load our Dataset-JSON files that we want to inspect.



One can either use drag-and-drop, or use the "Add" button, and then use the file chooser that is displayed. Once one or more files have been selected, also the "Remove" and "Clear" buttons become available. The "Remove" button is used to remove some files from the list after having selected them, and the "Clear" button to remove all files from the list.

Also here, the "Add from API" button is for demonstration purposes only and will be explained in the section "Using APIs".

When loading submission files, it is always a good idea to load the DM (Demographics) dataset in the case of SDTM and SEND, or the ADSL (ADaM Subject-level Analysis) in the case of ADaM, as this will not only to validate some of the data against DM/ADSL, but also to "toggle" between a record in any dataset and the "basic" subject data in the DM/ADSL dataset.

Let us first add a few datasets already, e.g.:



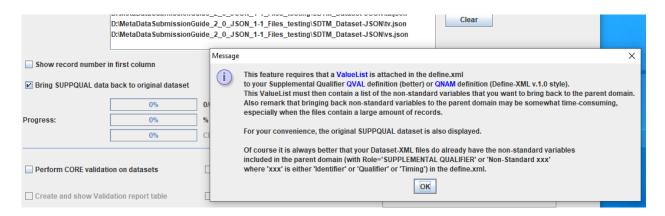
Don't mind about the order of the datasets in the list, the system will itself set up a more natural order.

Depending on which standard was selected, and whether one will use the define.xml for the metadata or not, a number of features become available. We explain them here for the case that one wants to use the define.xml for the metadata, as this should always be the preferred way.

		uide_2_0_JSON_1-1_Files_testing\SDTM_Dataset-JSO uide_2_0_JSON_1-1_Files_testing\SDTM_Dataset-JSO	•
Show record number	in first column		
Bring SUPPQUAL data	a back to original dataset		
	0%	0/0 files read	
Progress:	0%	% validation done	
	0%	CDISC Library	
Perform CORE validat	ion on datasets	☐ Create HTML report from extended JSON report	

First of all, when the checkbox "Show record number in first column" is checked, the system will add an extra column as the first column, and assign a record number to each record, which is fixed to that record. This also means that when later applying sorting, filtering, or any action that rearranges the order of the records, the record number will remain the same. Such a "pinned" record number can then enormously help when referring to a specific record later. Also see the section "sorting and filtering".

The second checkbox "Bring SUPPQUAL data back to original dataset" will only be enabled when the define.xml is to be used for the metadata. When one selects, some additional information is shown:

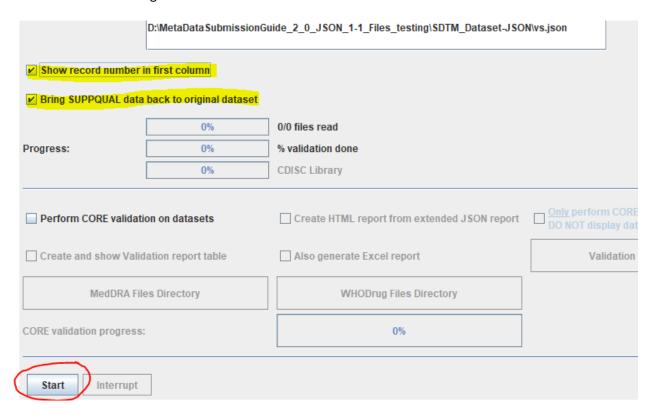


Essentially, it states that this feature will only work when the define.xml has "ValueLists" attached to QVAL for the variable definitions in the SUPPxx datasets.

As we have DM and SUPPDM, and EC and SUPPEC in our list, we check the checkbox. We will discuss the results a bit later.

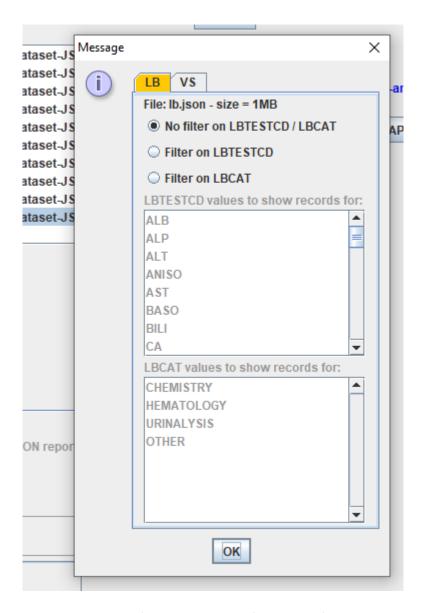
At this moment, we do not set any special options, using the "Options" button yet, so that only essential validation will be performed. We will later explain some of these in a separate section.

With the current settings:



we now click the "Start" button to start the viewing process.

First, the system loads the define.xml, which can take some time, depending on the complexity of the contents of the define.xml (progress bar), in most cases followed by the following dialog:



This dialog is only displayed when some of the to-be-loaded files have a file size that exceeds a certain threshold (see later when treating all the options) and the define.xml has codelists attached to variables such as "--TESTCD", "--CAT" allowing to load subsets of data, such as only "hematology" data for LB. This feature is expected to facilitate the reviewing process considerably.

As this dialog is based on information about codelists in the define.xml, it is skipped when the user has decided to use the metadata from the Dataset-JSON files themselves.

We will first not do any "pre-loading filtering".

After clicking "OK", the loading process proceeds and after a few seconds (or more, depending on the number of records to be loaded), the window with all the tables is displayed:



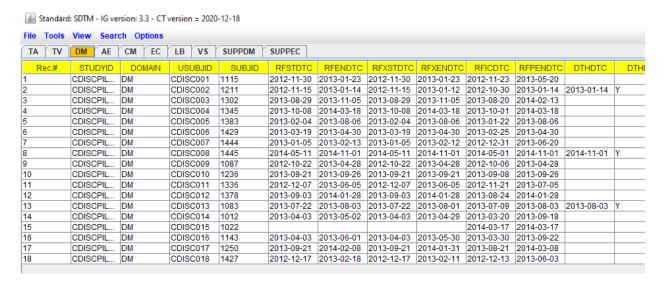
¹ When using the classic "SAS Universal Viewer", the user has no other choice than loading all the records, often leading to "pagination". When then doing additional filtering, the filter only applies to the current "page". As such, using the "SAS Universal Viewer" can easily lead to review problems.

with one tab for each dataset.

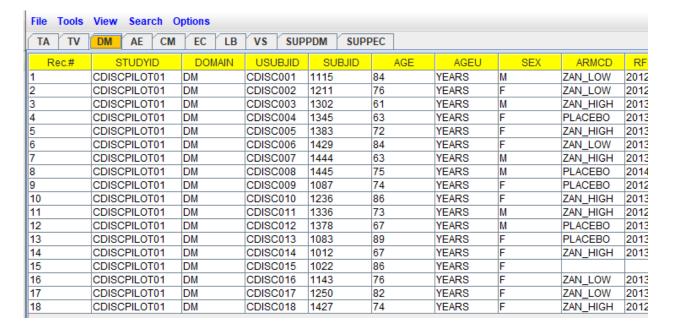
As one sees, the system has reordered the datasets into a more natural order, with the "Trial Design" datasets coming first, followed by DM (Demographics) and with the "Supplemental Qualifier" datasets coming at the end

Viewing the results - basic features

Let us first inspect the results for the DM dataset, by clicking on the DM tab:



As one sees, and extra column with the (generated) record number has been added One can now resize columns with the mouse, or even move columns from one place to another (i.e. changing the order of the columns) which is always very useful, as one can then put the columns of interest together. For example, when one is mainly interested in age, sex, and assigned arm of the subjects, one can drag the columns in the following order:



Viewing values for Supplemental Qualifiers

As we checked the checkbox "Bring SUPPQUAL data back to original dataset", we should now also see the values of the "non-standard variables", which were "banned" to SUPPDM, in the DM dataset itself. Remark again that this will only work when a "ValueList" was assigned to the QVAL variable in the define.xml.

When scrolling to the right, we find:

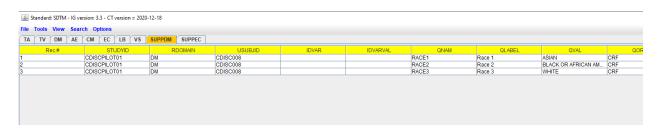
										I		
BRTHDTC	RACE	ETHNIC	ARM	ACTARMOD	ACTARM	ARMNRS	ACTARMUD	COUNTRY	RACE1	RACE2	RACE3	RA
1928	WHITE	NOT HISPA	Zanomalin	ZAN_LOW	Zanomalin			USA				
1936	WHITE	NOT HISPA	Zanomalin	ZAN_LOW	Zanomalin			USA				
1951	WHITE	NOT HISPA	Zanomalin	ZAN_HIGH	Zanomalin			USA				
1950	WHITE	NOT HISPA	Placebo	PLACEBO	Placebo			USA				
1941	WHITE	NOT HISPA	Zanomalin	ZAN_HIGH	Zanomalin			USA				
1929	WHITE	NOT HISPA	Zanomalin	ZAN_LOW	Zanomalin			USA				
1949	WHITE	HISPANIC	Zanomalin	ZAN_HIGH	Zanomalin			USA				
1939	MULTIPLE	NOT HISPA	Placebo	PLACEBO	Placebo			USA (ASIAN	BLACK OR AFRICAN AMERICAN	WHITE)	
1938	WHITE	NOT HISPA	Placebo	PLACEBO	Placebo			USA				
1927	WHITE	NOT HISPA	Zanomalin	ZAN_HIGH	Zanomalin			USA				
1939	WHITE	NOT HISPA	Zanomalin	ZAN_HIGH	Zanomalin			USA				
1946	BLACK OR	NOT HISPA	Placebo	PLACEBO	Placebo			USA				
1924	WHITE	NOT HISPA	Placebo	PLACEBO	Placebo			USA				
1945	WHITE	NOT HISPA	Zanomalin	ZAN_HIGH	Zanomalin			USA				
1928	WHITE	NOT HISPA				SCREEN F		USA				
1936	WHITE	NOT HISPA	Zanomalin	ZAN_LOW	Zanomalin			USA				
1931	WHITE	HISPANIC	Zanomalin	ZAN_LOW	Zanomalin			USA				
1938	BLACK OR	NOT HISPA	Zanomalin	ZAN_HIGH	Zanomalin			USA				

that there is one record having "MULTIPLE" for "RACE", and the individual races being depicted in "RACE1", "RACE2", and "RACE3". Also here, we can drag columns so that we get a better oversight of what we are really interested in:

SUBJID	AGE	AGEU	SEX	ARMCD	RACE	RACE1	RACE2	RACE3
1115	84	YEARS	M	ZAN_LOW	WHITE			
1211	76	YEARS	F	ZAN_LOW	WHITE			
1302	61	YEARS	M	ZAN_HIGH	WHITE			
1345	63	YEARS	F	PLACEBO	WHITE			
1383	72	YEARS	F	ZAN_HIGH	WHITE			
1429	84	YEARS	F	ZAN_LOW	WHITE			
1444	63	YEARS	М·	ZAN_HIGH	WHITE			
1445	75	YEARS (M	PLACEBO	MULTIPLE	ASIAN	BLACK OR AFRICAN AMERICAN	WHITE
1087	74	YEARS		PLACEBO	WHITE			
1236	86	YEARS	F	ZAN_HIGH	WHITE			
1336	73	YEARS	M	ZAN_HIGH	WHITE			
1378	67	YEARS	M	PLACEBO	BLACK OR			
1083	89	YEARS	F	PLACEBO	WHITE			
1012	67	YEARS	F	ZAN_HIGH	WHITE			
1022	86	YEARS	F		WHITE			
1143	76	YEARS	F	ZAN_LOW	WHITE			
1250	82	YEARS	F	ZAN LOW	WHITE			

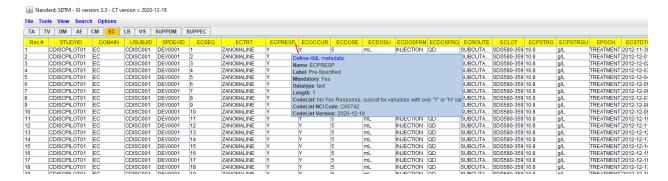
so that we now see that this subject is male, 75 years old, and received placebo.

Of course, we can still inspect the SUPPDM dataset itself by clicking the SUPPDM tab:



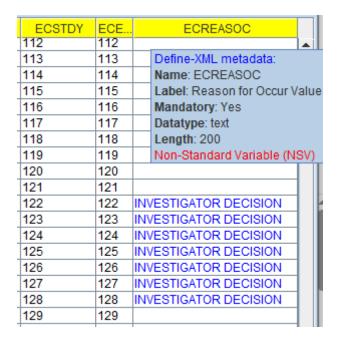
Showing the reviewer the same information, but in a much less user-friendly way.

Let us know at EC (Exposure as Collected) and SUPPEC. For EC we find:

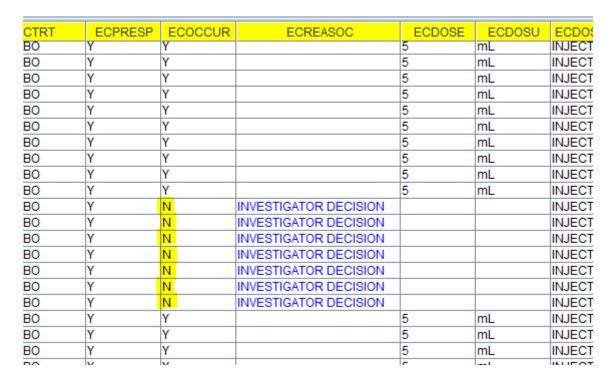


One sees that when hovering the mouse over a column header, one gets the metadata information as a tooltip.

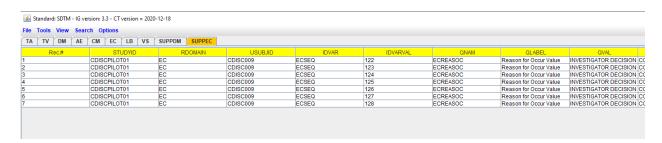
As we have a SUPPEC dataset, and checked the checkbox "Bring SUPPQUAL data back to original dataset", we also expect one or more additional columns for the NSVs (Non-Standard Variables). When we scroll to the right we indeed find:



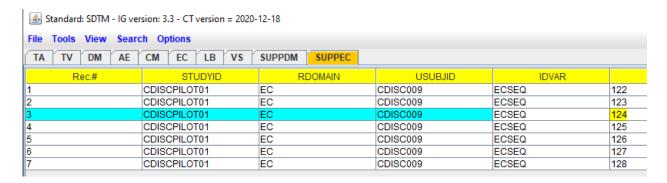
As this is related to non-treated subjects, we can better move that column, e.g.:



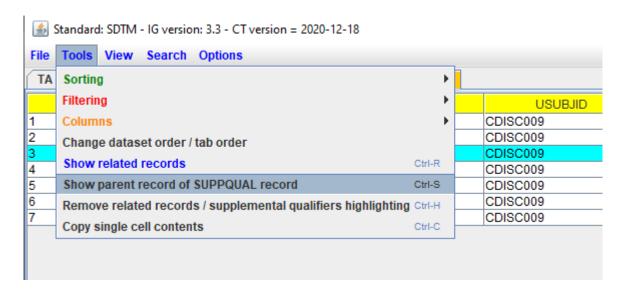
And the SUPPEC dataset:



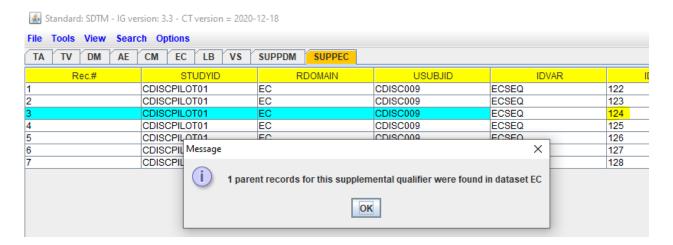
Even when one has <u>not</u> checked the checkbox "Bring SUPPQUAL data back to original dataset", one can easily find the parent record of a SUPPQUAL record. This can be done by selecting a row in the SUPPQUAL table (here SUPPEC):



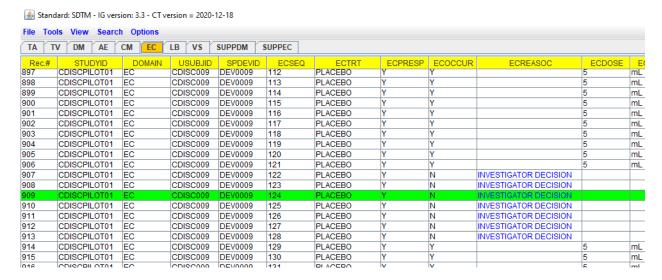
and then using the menu "Tools - Show parent record of SUPPQUAL record" (or just using Ctrl-S on the keyboard)



leading to:



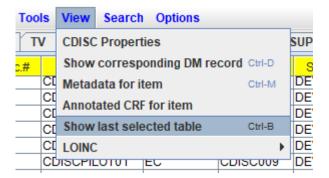
and after clicking "OK", the row for subject CDISC009 and ECSEQ=124 is selected (also the tab for EC is automatically selected) and highlighted:



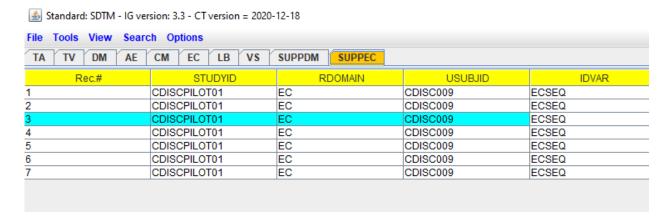
Remark that when the "identifying variable" (IDVAR) is e.g. a --CAT variable, this will usually lead to more than 1 parent record, and all these will then be highlighted.

One can now "toggle" between the record in EC and the record in SUPPEC (also when the checkbox "Bring

SUPPQUAL data back to original dataset" was <u>not</u> checked) by using the menu "View - Show last selected table", or easier and much faster by just using Ctrl-B on the keyboard:

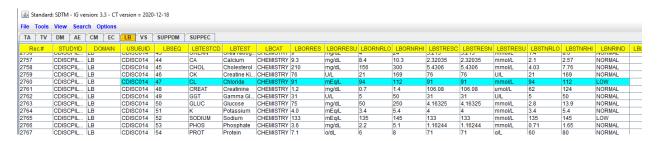


showing SUPPEC again with the selected record:



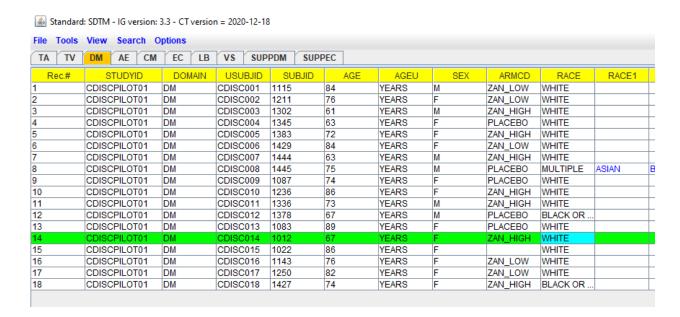
Looking up basic subject information in DM or ADSL

When inspecting a record in any of the datasets, one can always quickly "jump" to the corresponding record in DM (for SDTM and SEND) or ADSL (for ADaM). For example, when having selected a record in LB:



for subject CDISCPILOT14 with a low Chloride value, we can "jump" to the DM record for that subject by using the menu "View - Show corresponding DM record"², or much faster, just using Ctrl-D on the keyboard, leading to:

² In the case of ADaM, the menu item will show "Show corresponding ADSL record"



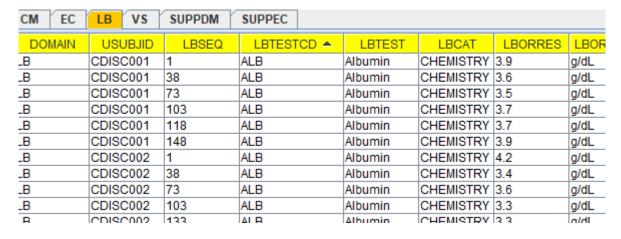
showing us that the subject was in the "high xanomeline" arm.

Also here, one can then use Ctrl-B to "toggle" between the two datasets and records.

Sorting and filtering

The result tables can easily be sorted and filtered,

One can sort within a column by clicking on the column header. For example, in LB by clicking on the LBTESTCD column header, the rows will be ordered in alphabetic order for the value of LBTESTCD:



The next click on the column header then sorts the rows in the reverse order:



Remark that for numeric variables ("integer" or "float" in the metadata) the rows will be sorted numerically for the variable value.

Sorting can also be accomplished using the menu "Tools - Sorting":



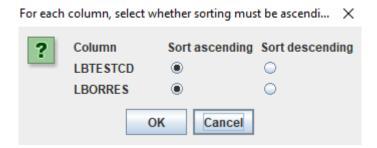
with keyboard shortcuts Ctrl-T (for sorting) Ctrl-U (for unsorting). When using "Sorting", this leads to:



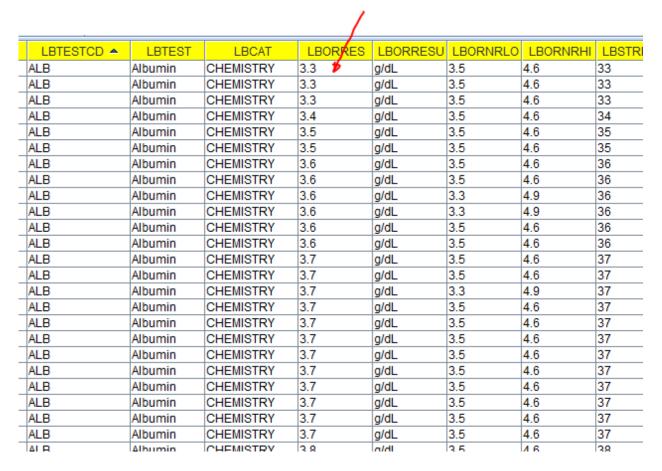
where one can add and remove columns to be sorted on. For example, for sorting on LBTESTCD and LBORRES:



This leads to a new dialog:

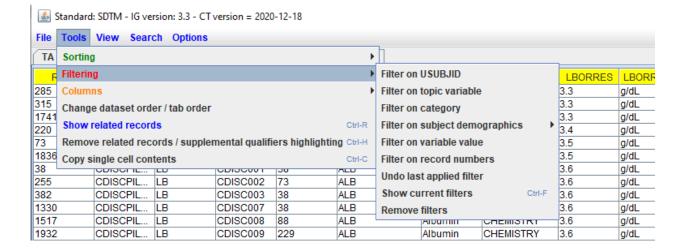


where one can decide to sort ascending or descending for each of the variables, e.g. leading to:



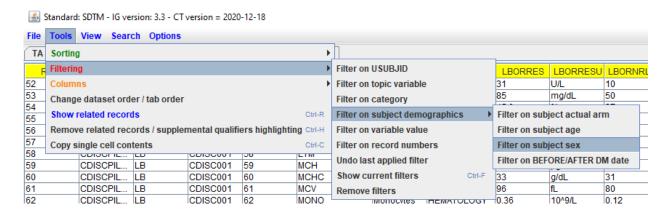
Bringing the table back to the original order can then be done using Ctrl-U.

Also filtering is easily possible, using the menu "Tools - Filtering":



One can filter on the subject, one the topic variable (--TESTCD for "Findings", --TERM for "Events", --TRT for "Interventions", on category (--CAT variable), on subject demographics (e.g. SEX, RACE, AGE, ...), but also on variable values (ranges) and of course on record numbers. One can apply "filters on filters", so making combination of filters, display them, and remove all the filters.

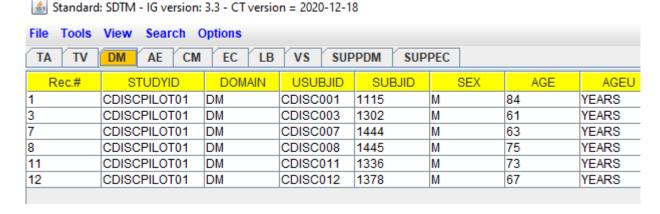
For example, if we want to filter on the "ketone" lab values for the male subjects, we start by applying a filter on "subject demographics", using the corresponding dialog:



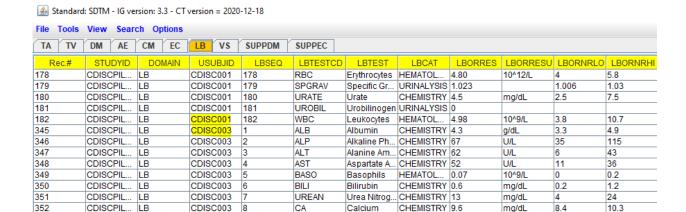
followed by:



and ensure that the checkbox "Apply to all tables" is checked. For DM, we will then see:



Only showing the male subjects: CDISC001, CDISC003, CDISC007, CDISC008, CDISC011 and CDISC012. In the LB table we then find:



CHOL

only containing records for male subjects.

LB

LB

CDISCPIL

CDISCPIL.

CDISCPIL..

353

354

On top of this, we want to apply a filter for LBTESTCD = KETONES.

CDISC003

CDISC003

CDISC003

8

We can do this by filtering on the "topic variable", as for LB, LBTESTCD is the "topic variable", or we can use "Filter on variable value", select LBTESTCD:

Calcium

CHEMISTRY 9.6

Cholesterol CHEMISTRY 255

Creatine Ki... CHEMISTRY 444

mg/dL

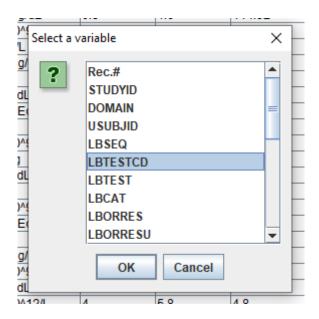
mg/dL

U/L

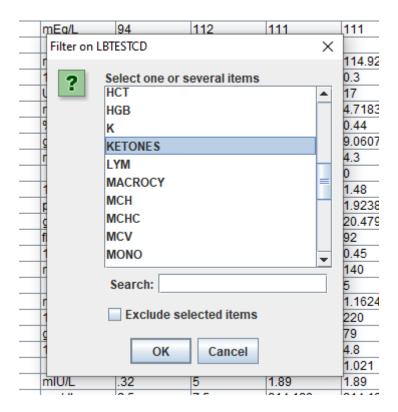
149

286

198

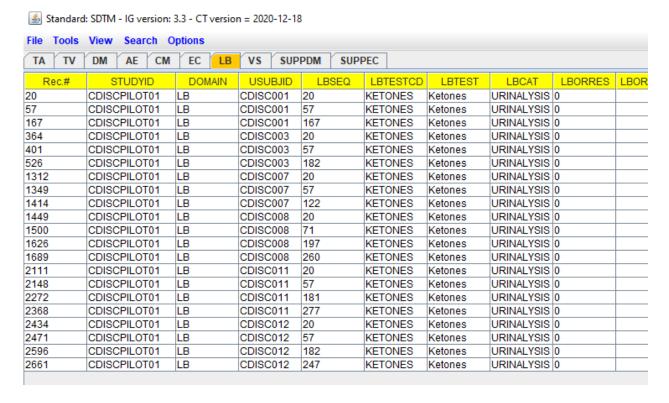


and then select "KETONES":



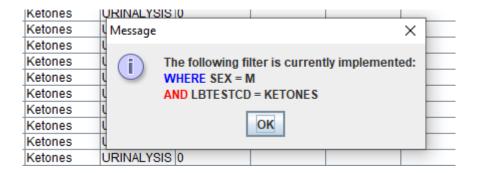
Remark that If one wants to select all values except for a few ones, one can also use the checkbox "Exclude the selected items".

The result after clicking "OK" then is:

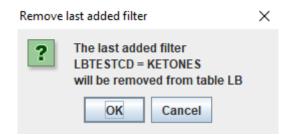


only showing the lab records for "ketones" for male subjects CDISC001, CDISC003, CDISC007 etc..

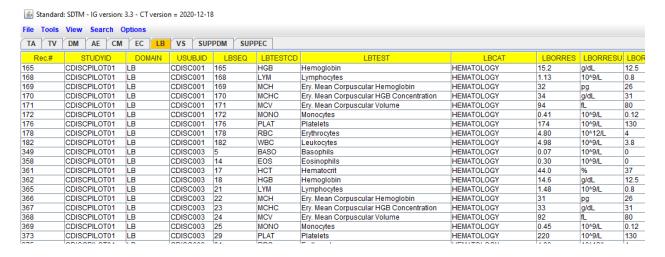
One can always check which filters are on, by using the menu "Tools - Filtering - Show current filters", in our case leading to:



One can also remove the last applied filter using the menu "Tools - Filtering - Undo last applied filter". A message is then displayed:



If we do so, and then add a new filter on LBCAT = HEMATOLOGY, the result is:

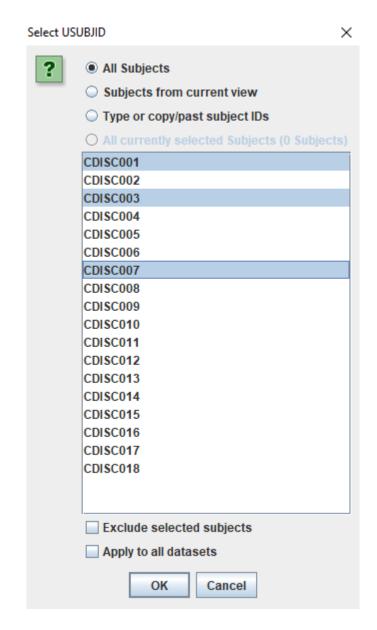


To remove all filters, use the menu "Tools - Filtering - Remove filters", leading to:



Usually, one will want to remove the filters on all the loaded datasets, otherwise things can get a bit messy. When then using "Yes, remove on all datasets", all the tables will be shown in their original state.

One can also always filter on specific subjects, by USUBJID, using the menu "Tools - Filtering - Filter on USUBJID". The system will then present a list of subjects by USUBJID, from which one can select the ones of interest, e.g.:



Also here, one can filter by exclusion, and apply the filter on either the current loaded datasets on all loaded datasets.

Remark that it is often to do "pre-loading filtering" as explains in the section "Filtering before loading".

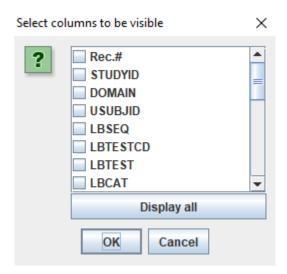
Moving columns around and hiding columns

We already showed how columns can be moved from one place to another, simply by dragging the column header. As in SDTM and SEND, the column order is in principal fixed, one can reset to the original order by using the menu "Tools - Columns - Reset column order":

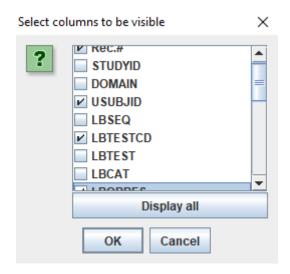


The menu "Tools - Columns - Move selected column" can also be used to move a column to another place (right-click on the column header first to select it), but usually just dragging the column will be the easier way.

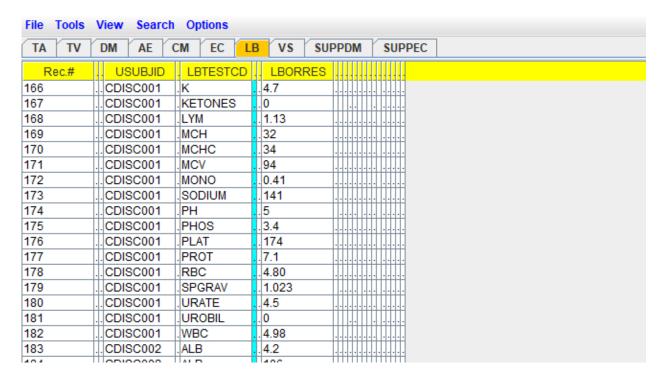
One can also hide (and later re-display) columns by using the menu "Tools - Columns - Hide/Display columns". This brings the following dialog up:



When we e.g. want to have "Rec.#", USUBJID, LBTESTCD, LBORRES and LBORRESU displayed, we select these:



leading to:



Remark that the columns to be hidden are not completely hidden, but just "minimized", so that dragging on the edge can make individual columns visible again.

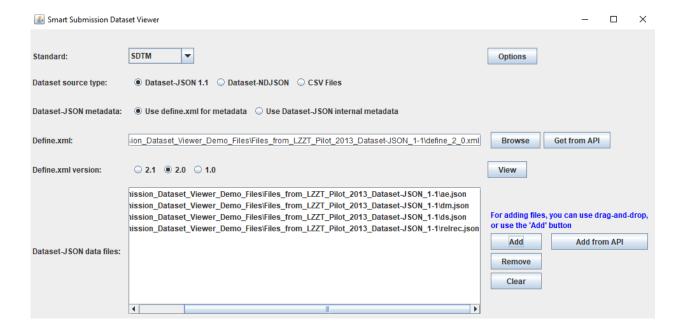
We might want to change this behavior in future depending on the user responses.

Making all the columns visible again, i.e. giving them their original width, can then be done using the same menu, and clicking "Display all".

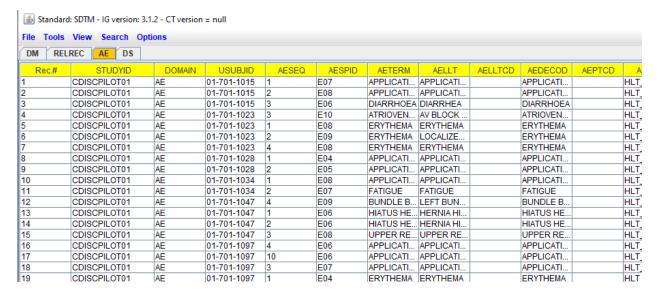
Connecting Related records (RELREC)

We did already show how one can establish the relationship between rows in the SUPPQUAL (supplemental qualifier) datasets and the "parent" records in the "parent" table, either by "bringing back" the "non-standard variables" to the parent dataset as additional columns (which however requires a quality define.xml) or by "toggling" between the related rows in both the datasets.

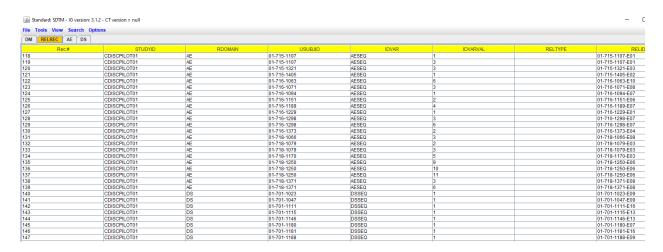
There is however another importantmechanism of having relationships between records and tables, by the use of the "RELREC" (Related Records) datasets. It is often used to establish the relationship between the AE (adverse events) dataset and the CM (concomitant medications) datasets. Unfortunately, the sample dataset from the "Metadata Submission Guidelines" does not have a good example, so we need to take another example. We find this in the files from the original LZZT pilot 2013 which we converted from XPT to Dataset-JSON 1.1, and which contain records for describing the relation between AE (adverse events) and DS (dispostion). Loading these, together with DM (always a good idea) into the software:



and start loading by clicking the "Start" button, this then leads to:



and for the RELREC table:



When sorting by RELID, we better see which AE and DS records are related:



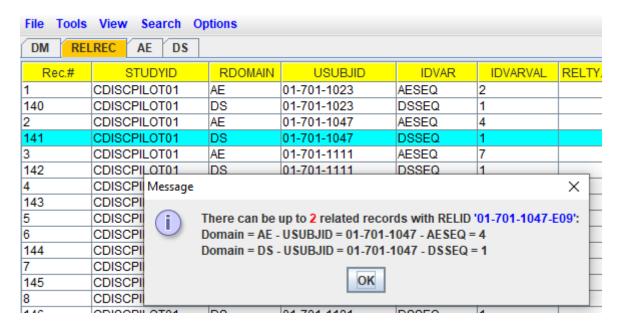
showing that for subject 01-701-1023, the AE record with AESEQ=2 is related to the DS record with DSSEQ=1,

and that for subject 01-701-1047, the AE record with AESEQ=4 is related to the DS record with DSSEQ=1.

If we select the latter (record # 141):



and then use the menu "Tools - Show related records" (or simply use Ctrl-R on the keyboard), this leads to a message:



and after clicking "OK", additional information is displayed in a separate window:

Related Records core properties and values

RELID: 01-701-1047-E09

Study ID	Dataset	Subject ID	AESEQ	AETERM	VISITNUM/ VISIT	AEDTC	AESTDTC	AEENDTC
CDISCPILOT01	AE	01-701-1047	4	BUNDLE BRANCH BLOCK LEFT		2013-03-10	2013-03-10	
Study ID	Dataset	Subject ID	DSSEQ	DSTERM	VISITNUM/ VISIT	DSDTC	DSSTDTC	DSENDTC
CDISCPILOT01	DS	01-701-1047	1	ADVERSE EVENT	6 (AMBUL ECG REMOVAL)	2013-03-29	2013-03-29	

showing the details of the relation.

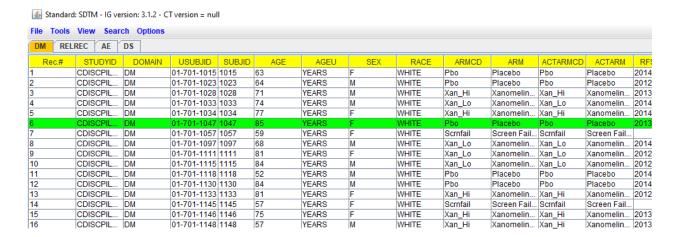
If we then select the AE tab, the specific record is already highlighted:



and when selecting the DS tab, also there, the specific record is already highlighted:



one can then "toggle" between the DS, AE, and RELREC datasets using Ctrl-B on the keyboard, and using Ctrl-B (or the menu "View - Show corresponding DM record"):

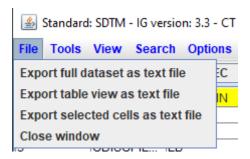


easily find out that this relation is about a 85 year old female who was in the placebo arm.

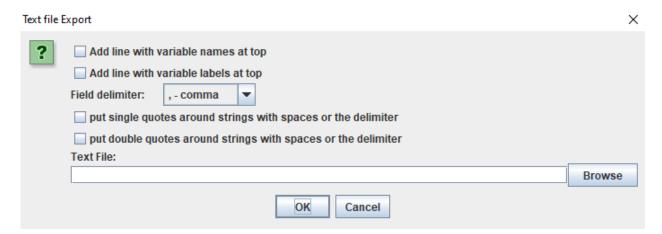
Exporting the table

In some cases, one want to export the table contents, of parts of it, as simple text or CSV files for use in other applications. Essentially, this should not be necessary, as we do already have the data in modern Dataset-JSON format, but as long as tools like Excel are used in clinical research (they shouldn't except for bookkeeping, that's what they were developed for), people will want to be able to import data from dedicated applications into such applications.

The "Smart Submission Dataset Viewer" provides several possibilities though the menu items "File - Export ...":

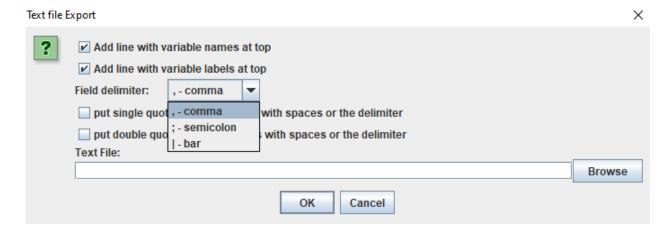


One can either export the full dataset of the selected table as a text file, or only the rows and columns that are currently visible, or only the set of currently selected cells. In each of these three cases, the system will ask <u>how</u> to export the data, by the following dialog:



The first two checkboxes allow to add a line with the variable names, and a line with the variable labels, which can serve as a header in receiving applications.

The system then asks what the delimited between the fields should be:



giving the choice between the comma, the semicolon and the vertical bar ("pipe" character).

The next two checkboxes allow to set whether either single or double quotes need to be added around the strings values of the variables. This can be important when e.g. the comma is selected as the field delimiter, but the string values themselves can contain the comma character³.

The text field and "Browse" button then allow to define where the text (or CSV) file needs to be written to. Remark that the file will be written using UTF-8 encoding, essentially representing Unicode.

An example of such an export is:

```
| New A. STUDYID, DOMAIN, USUNID, LENGE, LETTERCO, LETTERC, LEGGERS, LEGGER
```

Basic validation

During loading, one will have noticed two "progress bars", one showing the progress of the loading of the data, and the other one showing the progress of the basic validation that is always performed. This primary validation of the data is really very basic, e.g. checking whether "required" variables are populated. For example, when the column AEDECOD (Dictionary-Derived Term for the adverse event), which is a "required" variable, is not populated, the cells are colored red, and a tooltip is provided displaying the cause of the

³ The safest is always to use the vertical bar ("pipe" character) as the field delimiter, as it is a character that seldomly appears in the table values, but it might be that the receiving application, e.g. a worksheet, is not capable to accept that.

issue:

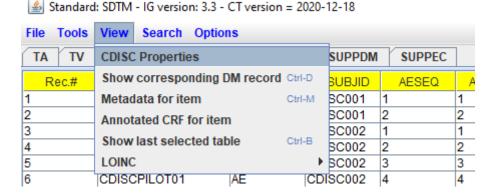
AELLT	AELLTCD	AEDECOD	AEPTCD	AEHLT	AEHLTCD	AEHLGT					
		ERROR: Required or expected variable value (AEDECOD)									

Further validation can be done by setting some options, by clicking the "Options" button in the main window, and by using the CDISC CORE (CDISC Open Rules Engine) features⁴ of the Smart Submission Dataset Viewer.

Search and View

In the "tables view", there are two other menus that we did not explain so far, but are easy to understand. When using the menu "View", one can choose between:

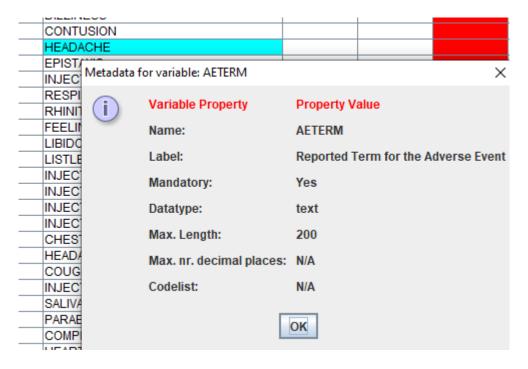
- CDISC properties:
- Show corresponding DM record (Ctrl-D): this was already explained before
- Metadata for item
- Annotated CRF for item
- Show last selected table (Ctrl-B): allows to "toggle" between tables
- LOINC => ... : this will be explained separately later



View - CDISC properties: displays which IG was used and which version(s), and which version(s) of the CDISC controlled terminology were used. This information is retrieved from the define.xml file.

View - Metadata for item: when a single cell has been selected, and this menu is used, a dialog with metadata information is displayed. For example when "headache" in the AETERM was selected, the resulting information is:

⁴ At the moment of writing (December 2024), the CORE engine is being adapted to accept Dataset-JSON 1.1 as being the input format. CORE was already able to work with Dataset-JSON 1.0, but as we are now moving to Dataset-JSON, some adaptions are necessary.

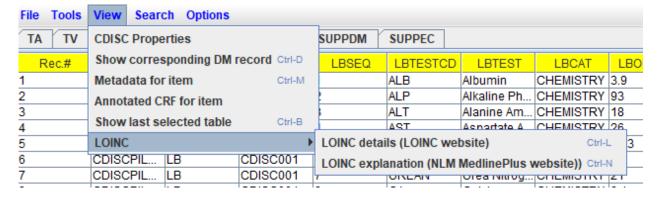


Also this information is retrieved from the define.xml

View - Annotated CRF for item: The goal of this feature is that when a single cell is selected, and the the link to the annotated CRF is present in the define.xml, the annotated CRF is opened in the default PDF viewer on the correct page, or, when there are several pages for the variable, on the first of these. This feature has not been implemented yet - it is planned for 2025.

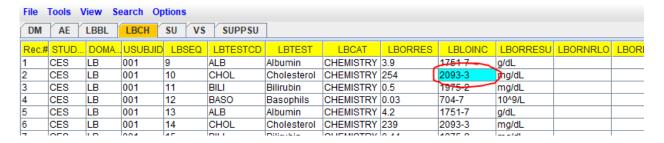
The next one "View - LOINC" is a bit special, but important due to the growing importance of LOINC as the only real identifier of the test, not only for lab tests, but also for vital signs, microbiology, virology, questionnaires, ECG, etc.⁵.

When selecting the menu "View - LOINC", we get additional possibilities:

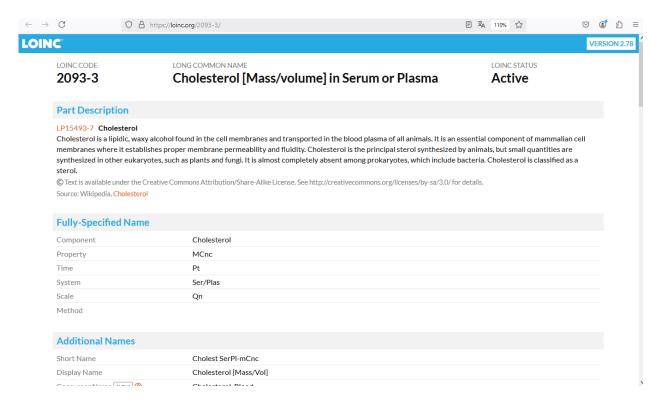


For example, when having selected the code "2093-3" in the LBLOINC column:

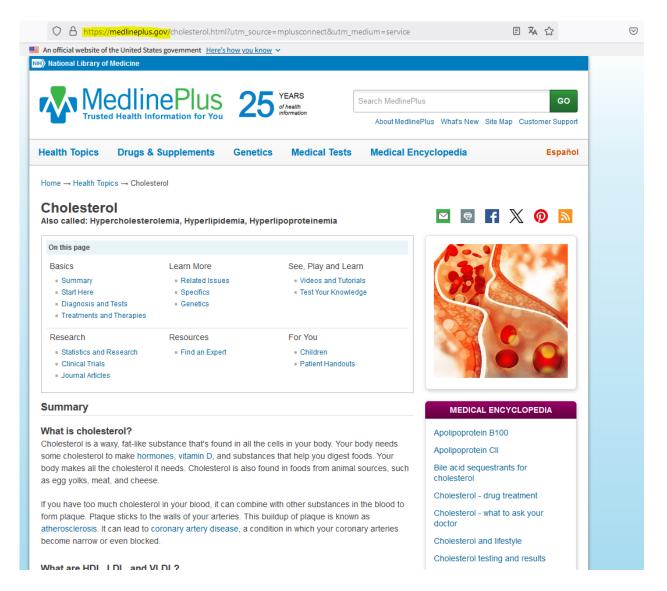
⁵ CDISC controlled terminology, as a "post-coordinated" system is not capable to uniquely identify tests, also not when combining values from variables such as LBTESTCD, LBSPEC, LBMETHOD, ... LOINC, as a "pre-coordinated" system makes this "uniqueness" possible. This is also the reason why LOINC is heavily used in healthcare, and is mandated by law in many countries (the number growing) for ordering lab tests and reporting lab results.



and using the menu "View - LOINC - LOINC details", this will trigger a process that opens the LOINC website page for this specific LOINC code in the user's default browser, i.e.:

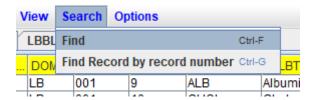


similarly, when the menu "View - LOINC - LOINC explanation" is used, the page for that code from the NLM (National Library of Medicine) is opened in the user's default browser, i.e.:

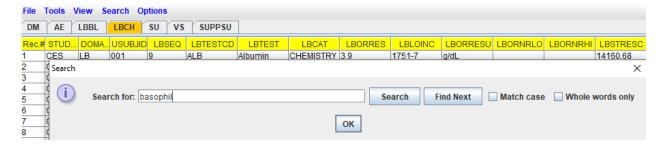


IMPORTANT REMARK: Essentially, this is an implementation of RWS (RESTful Web Services). The potential for such features is enormous! Using the mechanism, just by selecting the LOINC code, it would become possible to perform things such as literature searches, connect to specialist technical literature and other documentation, and also connect directly connect to specialist (in medicine or in clinical research) Al systems. Implementing this in the Smart Submission Dataset Viewer will often only require to add a few lines in the source code ... the only limitation being ... lack of creativity.

Another menu is "Search", allowing to quickly find a word, part of a word or value in the table, or to quickly jump to a specific record (by record number) in the selected table:



For the former (Search - Find), a separate dialog is then popping up:



which is pretty self-explaining.

Filtering before loading

When one has large datasets (thousand or even millions of records), it is very hard for reviewers to keep oversight. Of course one can filter after the dataset has been loaded, but in such cases, classic viewers such as the "SAS Universal Viewer" need to do "paging" in order to save memory. With "paging", any filters are only applied to a single page, which can easily lead to information loss. For example, when a reviewer only wants see the Albumin lab test results of a very large dataset, with the "SAS Universal Viewer" he/she can apply the filter on the first loaded page, and needs to repeat the whole process when wanting to see the Albumin values for subjects on other pages.

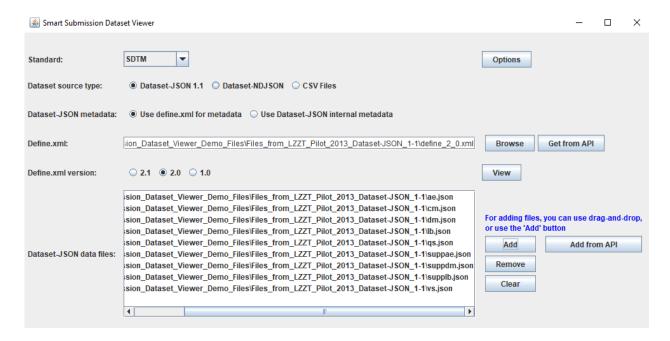
The Smart Submission Dataset Viewer uses a smarter method. Before loading, one can already apply a filter, and only those records that pass the filter are really loaded. This not only saves a lot of memory, but also ensures that <u>all</u> the records that pass the filter will be displayed.

In the Smart Submission Dataset Viewer, filtering is usually based on the --TESTCD or --CAT ("test code" and "category") variables for "Findings" datasets. In the future, we also want to apply the same principle for "Events" datasets (e.g. based on --DECOD), and for "Intervention" datasets, e.g. based on --TRT or --CAT or --INDC. We need input from the users about which variables are the most useful ones to do pre-filtering for these classes.

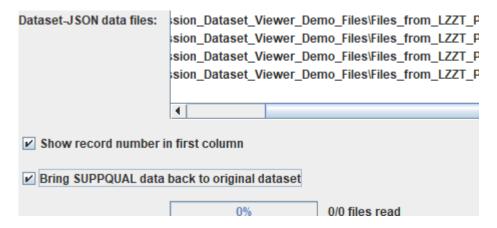
To find out which e.g. test codes or categories are available, a define.xml must be present and a codelist must be attached to the given variable. This is usually already the case for the "Findings" datasets.

The system always offer the user to do pre-filtering on specific datasets, when there is a codelist available from the define.xml, and <u>urge</u> the user to do pre-filtering based on the file sizes, the default being 20MB (this threshold can be changed by the user, see later). For example, in our case, the LB.json dataset has a size of 0.7MB (which is much less than for the SAS-XPT corresponding dataset, which has 2.7MB), but still has almost 3500 records. The LB dataset from the LZZT-pilot is 10.4MB in size (the corresponding LB.xpt is 66MB), and has almost 60,000 records, which means that in such a case, pre-filtering may very well make sense, as one cannot imagine how a reviewer can handle such an amount of records just by visual inspection.

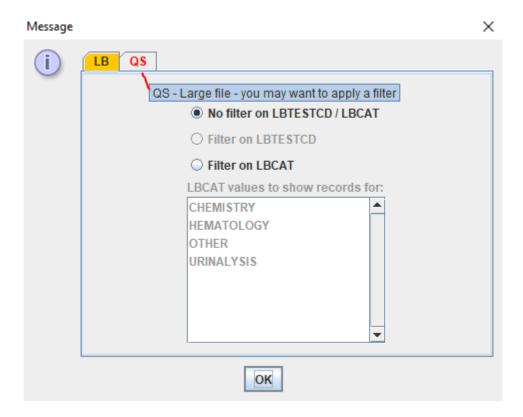
For demonstrating the pre-filtering, we take the SDTM Dataset-JSON files of the LZZT pilot, and e.g. select:



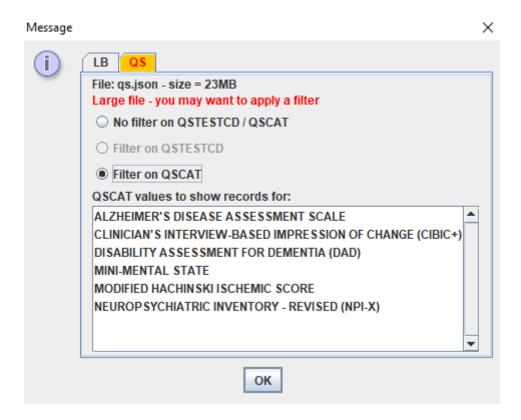
the QS dataset having the largest file size (23MB) followed by the LB dataset (10MB). We also load the SUPPLB dataset (8MB) and the SUPPAE and SUPPDM datasets, and ask the system to "Bring SUPPQUAL data back":



When then clicking "Start", first the define.xml is loaded, and then the system shows us the following dialog:



The system found codelists for LBCAT in the define.xml, so offers to filter on it, and strongly suggests to do pre-filtering for QS, for which it presents:



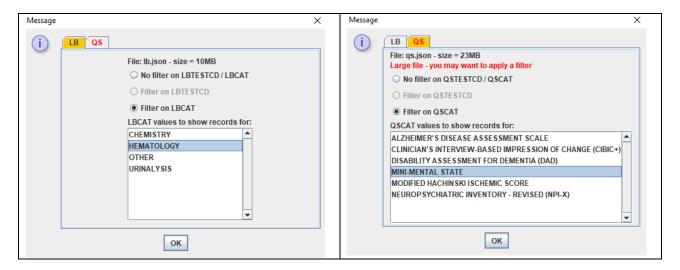
as it found a codelist for QSCAT, but none for QSTESTCD.

When we leave everything "as is", i.e. do not any pre-filtering, it takes 18 minutes to load, especially "recombining" SUPPLB with LB taking a lot of time.

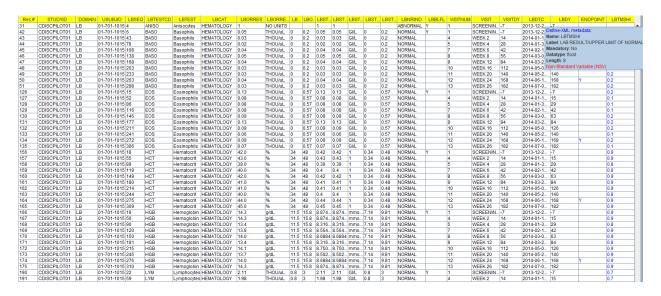
Anyway, visually inspecting 60,000 LB records doesn't make sense anyway isn't it?

So, let us be smarter, and do pre-filtering.

For LB, we pre-filter on "hematology" only, and for QS on "Mini-mental state":

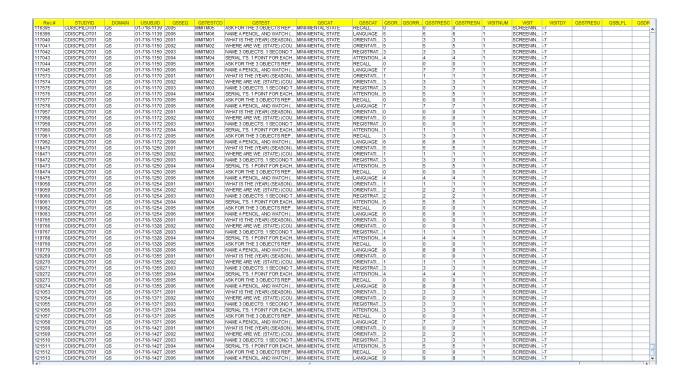


and start loading again. Loading the selection now takes 7 minutes, with the result for LB being:



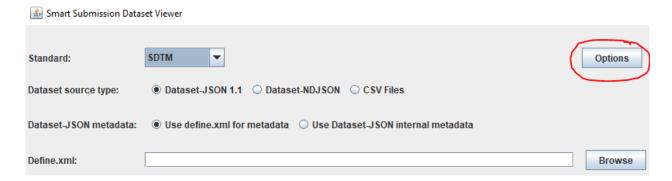
Remark that the record numbers are the same as if all the records were loaded.

For QS, we find:

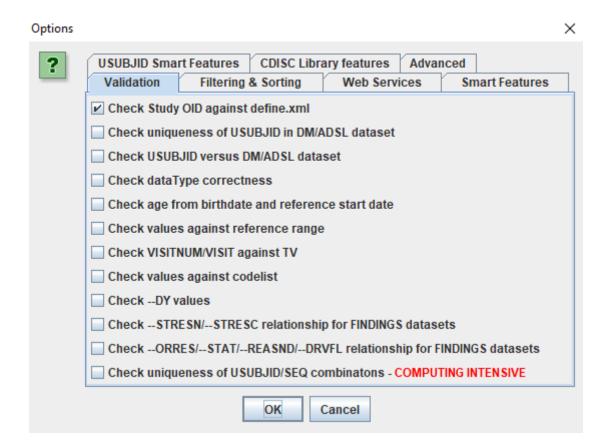


Setting Viewer Options

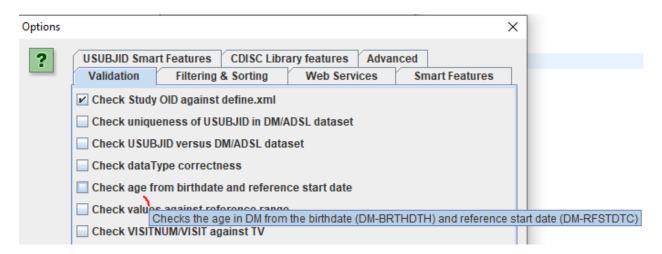
In the main window, one also sees an "Options" button on the right-top side:



When clicked, the following window with different tabs is shown:

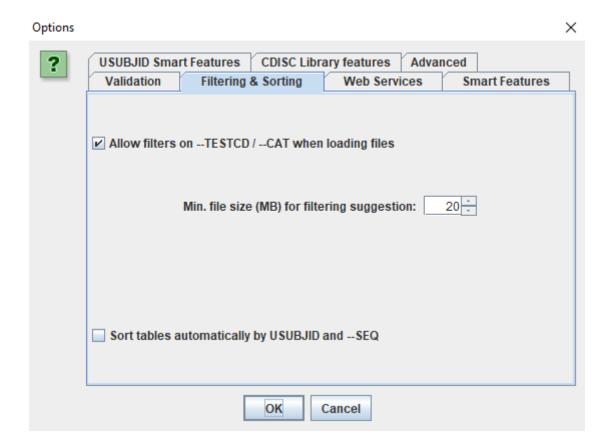


The options for the first tab "Validation" allows to switch on/off certain basic validation options. Most of them are obvious, but one can always hold the mouse over each of them to get more information. For example:



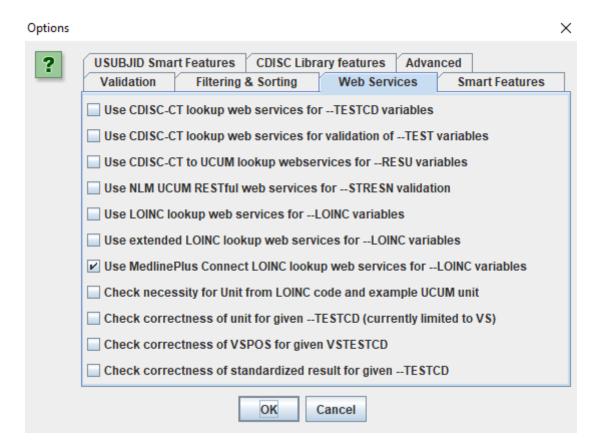
We will later, in the section "Validation" show how these can be used.

A second important tab is the "Filtering and Sorting" tab:

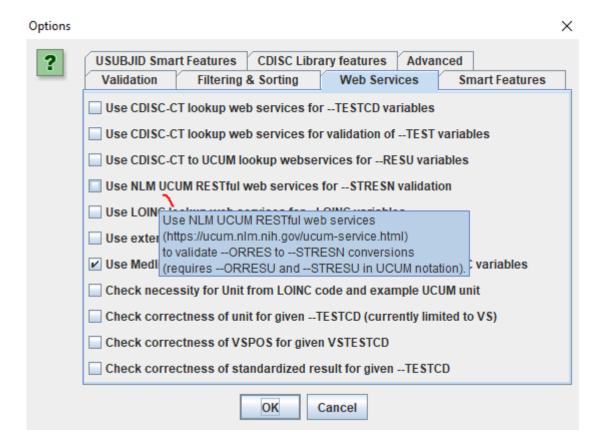


It allows the user to specify whether the system will propose to add "pre-loading filters" and if so, from what dataset size, the system will <u>urge</u> the user to apply a filter - see the section "Filtering before loading".

The next tab is about "Web Services":



The Smart Submission Dataset Viewer is unique in the ability to use (RESTful) web services to aid in validation. Also here, hovering the mouse over an item shows more details about the feature, e.g.:

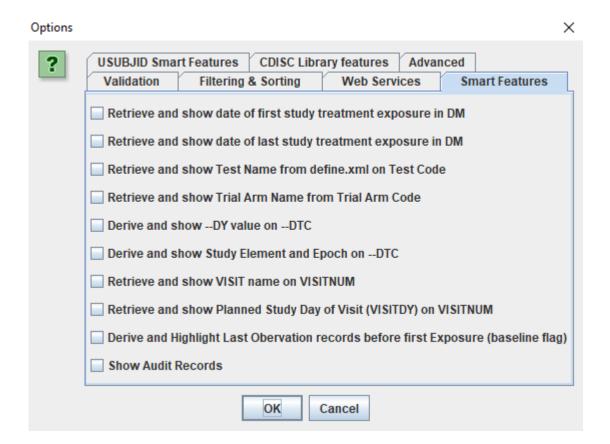


This feature allowing to use the NLM (National Library of Medicine) RESTful web services for automated validation whether the conversion the --ORRES values to the --STRESN values have been done correctly.

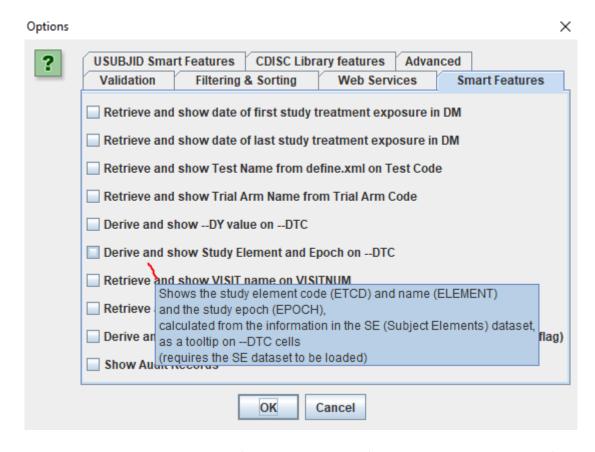
Some of these web service validation features even go far beyond what is available in CDISC CORE (see later) and especially what is offered by Pinnacle21 validation.

More explanation can be found in the separate document "<u>Using RESTful web services (NLM, FDA, XML4Pharma)</u>". Also this document will be updated in the near future.

The next tab is "Smart Features":



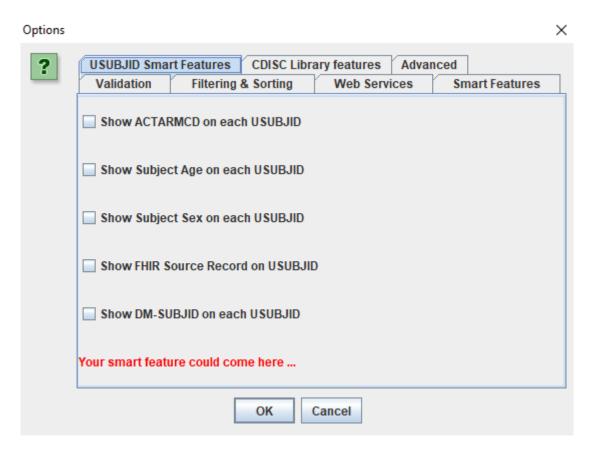
Also these are pretty self-explaining, with the possibility of obtaining additional information by howevering the mouse over an item, for example:



which, when checked, will show the value of EDTC and ELEMENT (study element code and name) as a

tooltip on each --DTC cell in each table, which can be very useful⁶. As one can expect, this requires the SE (Subject Elements) to have been loaded.

Furthermore, there is the tab "USUBJID Smart Features", which is mostly about getting information from the DM dataset (when loaded) in the tables for other datasets:



For example, when the checkboxes "Show Subject Sex on each USUBJID" and "Show ACTARMCD on each USUBJID" are checked, and one hovers the mouse over a "USUBJID" cell e.g. in the LB dataset, one obtains the information about sex and actual arm of that subject in the study. For example:

LB	01-716-1103	46	CHOL	Chole
LB	01-716-1103	81	CHOL	Chole
LB	01-716-1103	111	CHOL	Chole
LB	0 1-7 10-1 103	01-716-11	Chole	
LB	0 1-7 10-1 103	ACTARMO	Chole	
LB	01-716-1103	SEX: Male	Chole	
IR	01-716-1103	236	CHOL	Chole

immediately showing that the subject is a male person and the actual arm was "Xanomelin low". The same can of course be obtained by "toggling" between the current table and the DM table (using Ctrl-D and Ctrl-B), but just using hovering is of course faster.

As the Smart Submission Dataset Viewer is Open Source software (member of <u>COSA</u> - CDISC Open Source Alliance), additional "smart features" can always be added⁷ by users, the only limitation being lack of

⁶ This shows again that many "derived" variables in SDTM should not be in SDTM at all. They are only there as the reviewers at the regulatory authorities do not have access to modern viewers that understand CDISC submission standards. So they asked CDISC to add them in SDTM itself.

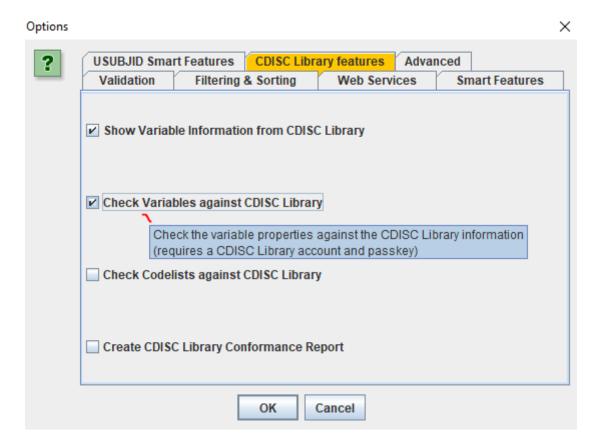
The main reason for this is the mandated use of SAS-XPT, and the use of primitive tools such as the "SAS Universal Viewer". Essentially, the Smart Submission Dataset Viewer lifts these limitations, so that in future, many of the "derived" variables could be removed from the SDTM standard again.

⁷ Of course, suggestions for new "smart features" are always welcome.

imagination ...

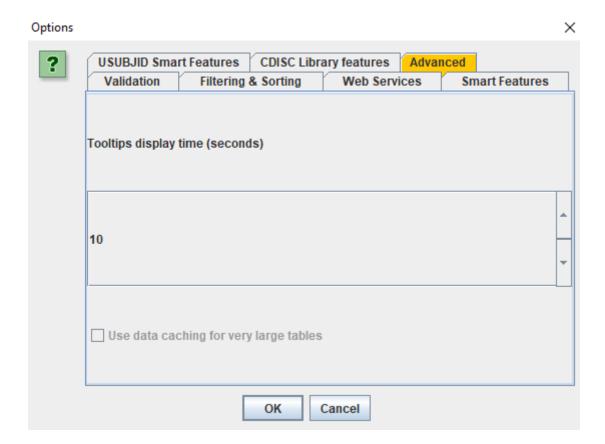
Very special is the the option "Show FHIR Source Record on USUBJID". This is an experimental feature showing how SDTM in Dataset-JSON format can include source records as e.g. from Electronic Health Records", and can then be visualized by a viewer. Currently, this has only be implemented for FHIR records. See the document "Support for embedded EHRs (FHIR)" for more details. The document still describes this feature for Dataset-XML as the transport format, but the same is of course true for JSON, as the FHIR standard has full support for JSON as the transport format. Remark that embedding source records in SDTM records is impossible when using SAS-XPT. Another reason to get rid of it as soon as possible.

Also the "CDISC Library" and its API have been implemented, tab "CDISC Library Features":



allowing to do validation by checking variable and codelist properties (from the define.xml) against the CDISC Library. This feature of course not only requires an internet connect. but also a CDISC Library account and API key. The latter must be stored in the "properties.dat" file - see section "The 'properties.dat' file".

The last tab is "Advanced":



For example, normally tooltips on cells or other items show up for 10 seconds by default, but the user can set this to another time duration.

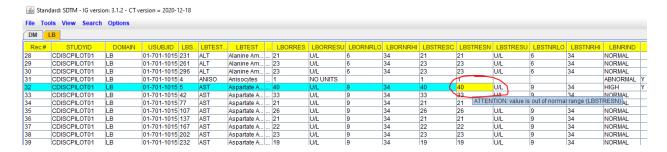
The option "Use data caching for very large tables" is currently still in development.

Extended validation

When none of the options in the "Validation" tab has been selected, only very basic validation will be performed (e.g. on "required" variables in the case of SDTM and SEND).

By switching on checkboxes in the "Validation" tab, additional validation checks are performed. This of course comes at the cost of loading time.

For example, when the checkbox "Check values against reference range" is checked, values that fall outside the --STNRLO / --STNRHI range (in any of the SDTM/SEND dataset) are highlighted. For example:



where the value of LBSTRESN for record 32, being 40 falls outside the range 9-34 U/L.

This feature can be used by reviewers to quickly find "out-of-range" values without needing to rely on the -- NRIND (Reference Range Indicator), provided by the sponsor, and which can have been incorrectly

assigned⁸. This means that the reviewer is now easily able to have an independent assignment, independent of what is provided by the sponsor. Many of the validation features of the viewer do allow the same independency.

The same e.g. applies to the feature "Check --DY values" (--DY = "Study day"). It is one of this variables where we see that very often errors are made, which pass unseen when e.g. using the "SAS Universal Viewer".

Another example is the use of the "CDISC Library". When in that tab, the checkbox "Check Variables against CDISC Library" is checked, and when hovers the mouse over a column header, additional information is provided, such as:

LBORRES	LBORRE	SU	LBORNRLO	LBORNRHI	LBS	STRESC	LBSTRES
21	U/L		6	34	21		21
23	U/L	De	fine-XML meta	adata:			23
23	U/L	Na	me: LBORRE	SU			23
1	NO UNITS	Label: Original Units					1
40	U/L	Mandatory: No					40
33	U/L	Da	tatype: text				33
21	U/L	Lei	ngth: 8				21
26	U/L	CD	ISC Library in	formation:			26
21	U/L	Na	me: LBORRE	SU			21
22	U/L	Lal	bel: Original U	nits			22
23	U/L	Co	re: Exp				23
19	U/L	ΧP	T simple data	type: Char			19
23	U/L	Co	deList NCLCo	de: C71620			23
19	U/L	Dis	crepancies w	ith CDISC Lib	rary:		19
0.05	THOU/uL	- N	o CDISC/NCI	codelist assig	ned	1	0.05
0.03	THOU/uL	Exp	ected: C7162	20			0.03
0.02	THOU/uL		0	0.2	0.02		0.02

also stating that it was expected that the CDISC codelist C71620 is assigned to variable LBORRESU, which isn't in this case.

We will not explain every validation feature further here as there are a lot of them, and new ones are regularly added. Important to know is that, except for the very basic ones, the default is that they are switched off.

Validation using CDISC CORE

Also CDISC CORE (CDISC Open Rules Engine) have been implemented in the "Smart Submission Dataset Viewer". At the moment of writing (December 2024), the feature has been disabled, as CORE is currently being updated for Dataset-JSON 1.1 (version 1.0 is now deprecated and should not be used anymore). As soon as v.1.1 of Dataset-JSON is implemented in CORE, we will re-add the features to use CDISC CORE validation.

A separate document explaining all the features of validation using CDISC CORE, as implemented in the Smart Submission Dataset Viewer, is currently in development.

⁸ This shows again that such variables such as --NRIND (Reference Range Indicator) should not be part of SDTM (as they are derived). Such variables have only be added on request by regulatory authorities, as the tools of the reviewers are not "smart" enough to derive them theirselves.

Using APIs

IMPORTANT REMARK: everything described in this section is purely experimental!

It is meant for demonstration purposes to show how APIs and RESTful Web Services (RWS) can be used to retrieve information from other or remote systems, not only full datasets, but also selected information, making the review process considerably more easy.

It also demonstrates that with the use of Dataset-JSON, it becomes in principle possible to evolve to "rolling submissions", in which SDTM/SEND and possibly also ADaM data already becomes available as soon as the first data points have been collected, as also the process of conversion from source data to especially SDTM and SEND can easily be made "rolling", thanks to modern conversion tools.

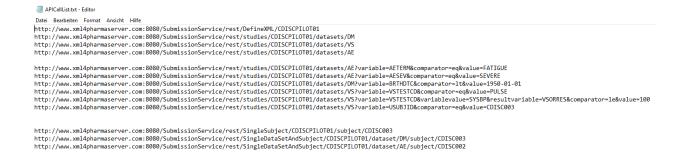
Such "rolling submissions" also would enable regulatory authorities to make final decisions within a short time after the last data point was collected, which would mean that new therapies can become available at least 1 to 2 years earlier than is currently the case.

In the main window, one may have noticed the buttons "Get from API" for define.xml, and "Add from API" for the datasets themselves:

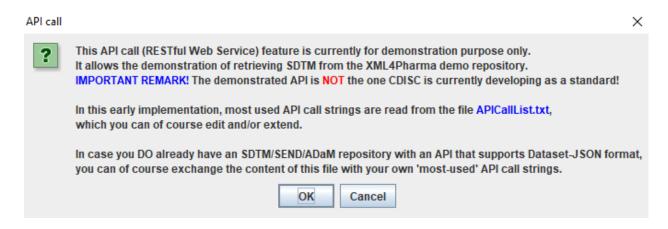


Just for demonstration purposes, we have build an SDTM dataset containing the information of the "LZZT Pilot" (no, not in XPT format ...), and developed a simple API, again, just for demo purposes, so nothing formal, to retrieve information from that database and load it into the viewer.

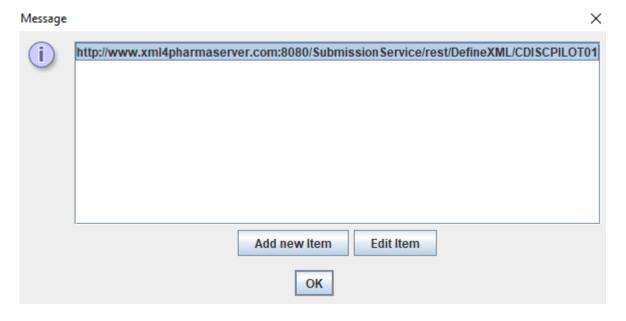
A number of "GET" queries have been added to a file "APICallList.txt", which comes with the distribution, and which can be edited/extended by the user:



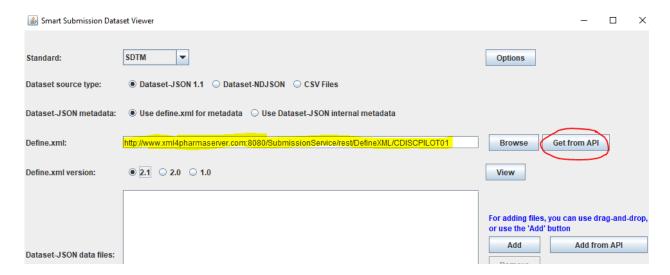
When then clicking the button "Get from API" for the define.xml, first some explanation is displayed:



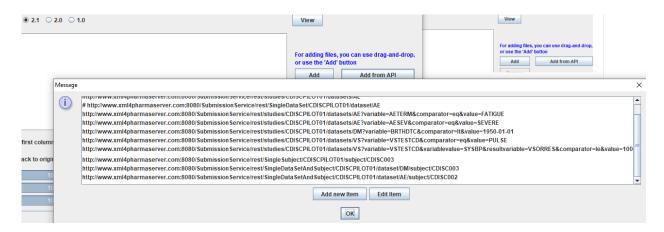
and then this list is read, and the available "GET" commands regarding the define.xml is displayed:



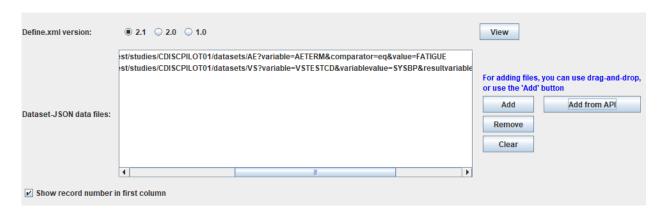
to which one can add new ones, or edit the existing ones, though this currently doesn't make sense, or it must be that one has implemented a similar API on an own server that has an SDTM database. When selecting the first one, and clicking "OK", the text field for define.xml is populated:



Similarly, when one clicks the "Add from API" button for the datasets themselves, the list is read and the "GET" instructions made available:



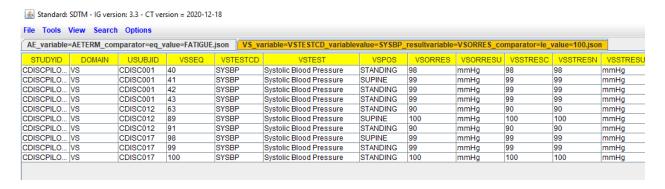
Just for the demo, let us select the one about "AE with AETERM is fatigue" and the one "VSTESTCD where VSTESTCD=SYSBP and VSORRES less or equal 100", leading to:



When then hitting "Start", the query is submitted to the server, and some Dataset-JSON is generated by the server and send back to our application as the "response", and then used by the viewer, leading to:



showing that the database has only 1 record with AETERM=FATIGUE, and



providing all the vital signs record for which the systolic blood pressure is less or equal to 100 mmHg.

Using such RESTful Web Services is much more efficient than having to load full datasets (or pre-filtered ones) and then having the reviewer to have to start filtering. It also does not require a high quality define.xml as is the case for the "pre-loading filtering" feature, which already can save a lot of time. One can easily imagine a simple application on top of the viewer that has a user-friendly graphical user interface and takes care of assembling the queries. And of course, one can also imagine Al to be used to do so, where the user simply talks to and says (or types): "Viewer, give me vital signs records for all the subjects with a systolic blood pressure lower than 100 for study XYZ ...".

Once again, please notice that what is presented in this section is not part of any formal standard, is not endorsed by CDISC, but was just developed to demonstrate the power of RESTful Web Services and API for reviewing submissions.

The "properties.dat" file

The "properties.dat" file allows to steer the behavior of the software. This file is immediately read after starting the software. The software comes with a sample properties.dat file in which some settings are "commented out". These are lines that start with a "#" character.

For example:

```
properties.dat - Editor
```

Datei Bearbeiten Format Ansicht Hilfe

#logfilepath=C:\temp

loglevel=INFO

- # Minimum amount of memory (in MB) that must always be available during loading #minmemorynecessary=5
- # Fill in your CDISC Library API key here if you want the CDISC Library features cdisclibraryapikey=abcdefghijklmnopqrstuvwxyz12345
- # Whether the question about converting XPT files to Dataset-JSON format is asked at startup skipXPTQuestion=false

In the first line (here commented out), the user can set to which folder the log files will be written. The default is that they are written to a folder named "logs" that is under the folder to which the software was installed.

The second line allows to set the "logging level". Under normal circumstances, one will set it to "INFO". If one encounters problems, best it to set it to "DEBUG". In that case, much more information is written into the log file, which you can then send to us to understand the issue, and helps the developers to further improve the software. Other possible settings are "WARNING" and "ERROR". For both of these, the log file will be much smaller, as it will only be written when a warning or error is thrown.

The next line is something that one mostly will have commented out. It defines the minimum amount of memory that must always be free during loading of the dataset. If, during loading, the amount of memory threatens to fall under this minimum, the system will issue a warning and try to do some "garbage collection". Change this setting only when you are a developer and have access to the source code. The "Smart Submission Dataset Viewer" has some validation options that use the "CDISC Library". In order to be able to use these, you will require a "CDISC Library API Key" which you can obtain from CDISC for free. See the CDISC Library API documentation for more details. When you have obtained one, copy-paste it into the "properties.dat" file.

When you start the software, you will be asked whether you first want to convert SAS-XPT files to Dataset-JSON, or whether you already have such, and wants to start with these.

When setting skipXPTQuestion=true", the system will skip this first dialog and immediately jump to the main window.