

# Much ADaM about nothing — a proc away in a day

**Endri, Rowland Hale**

i3, Berlin, Germany

CDISC is rapidly becoming adopted as the data standard for clinical trials and submissions of clinical trial data to the FDA. Within CDISC, ADaM datasets are an integral part of clinical study analysis and require significant data derivation to fulfill the needs of TLF provision. Can the creation of ADaM datasets be automated? Yes! This paper proposes an Excel-driven solution which greatly improves the efficiency and accuracy of ADaM dataset creation by defining ADaM dataset structure, complex variable derivation, and data checks within Excel. Using a library of SAS macros, the definitions then drive the automated production of SAS scripts which produce analysis ready datasets that meet the ADaM specification and which can be validated in accordance with regulatory requirements. Additionally, data check reports are created for Data Management to speed up the data cleaning process.

**Keywords:** CDISC, ADaM, Standard script, Automated mapping, Excel driven, Increased efficiency, SAS macros

## Introduction

CDISC ADaM version 2.1 is now released on the CDISC website. This document describes the fundamental principles of the ADaM analysis dataset and the design and purpose of submitted analysis datasets. The main principle of ADaM is ‘analysis-ready’ so that only minimal programming effort is needed to achieve statistical results from the ADaM datasets. In other words, although more than one statistical procedure may be required to calculate the statistical outputs, the goal of ADaM is to minimize programming effort and maximize concentration on the results.

The ADaM Implementation Guide version 1.0, which is posted on the CDISC website, describes in detail the standard structure of ADaM datasets, the variables they contain, and the derivation methods used. Using metadata to define the structure and derivation methods of ADaM datasets greatly facilitates the automation of ADaM dataset creation directly from CDISC SDTM datasets. This paper outlines the fundamentals of our approach and describes a step by step ADaM dataset automation process.

## Methods

### *ADaM datasets — principles of derivation*

Fig. 1 shows the flow of standard clinical trial data.

As CDISC SDTM, which contains only a small number of derived variables, becomes the standard structure for clinical trial data, ADaM datasets contain all derived variables required for the analysis. This

means that ADaM datasets should aim to be analysis ready, or ‘one proc away’ from the statistical outputs.

The ADaM Implementation Guide describes two standard data structures:

1. Subject-Level Analysis Dataset (ADSL Dataset). This dataset is the minimum requirement for ADaM. It provides key information for each subject in the clinical trial and contains a single record per subject;
2. Basic Data Structure (BDS Dataset). The primary keys of the BDS dataset are subject, analysis parameter and (depending upon the analysis) analysis timepoint.

Although some variables are required to exist within an ADaM dataset, e.g. AVAL, PARAM, DTYPE, etc., BDS is flexible in terms of additional rows and columns. This allows BDS datasets to provide robust and flexible support for most types of derivation and statistical analysis. Further information about these ADaM variables can be found in the ADaM Implementation Guide version 1.0.<sup>1</sup>

Understanding the structure of CDISC SDTM Datasets and the ADaM Implementation Guide helps us to automate the derivation of ADaM Datasets.<sup>2</sup> In general, the following steps are required to create ADaM datasets (See Fig. 2).

### **ADaM Metadata and Derivation Methods for the Automation Process**

As described above, all variables for each ADaM dataset, their attributes, and their derivation method are specified in the metadata. For clarity we define each analysis dataset in a separate Excel worksheet.

Correspondence to: Endri, i3, Joachimstaler Strasse, 12, 10719 Berlin, Germany. Email: endri.endri@i3global.com



Figure 1 Flow of standard clinical data

This makes it easier to add new variables to a particular dataset if needed. Keeping the metadata as simple as possible without compromising flexibility is the key to implementing the automation process efficiently and successfully.

The analysis variables derived for one particular ADaM dataset, e.g. ADEG, might be different from those derived for another, e.g. ADVS, and all of these should be clearly documented and linked in the metadata to ensure traceability. Besides information about the content, structure, and source of the data, the metadata also indicates the derivation or imputation method of the analysis variable.

The analysis dataset metadata has the following configuration and key variables:

- spreadsheet name: used to identify the analysis dataset name;
- variable name: variable name;
- variable label: variable label;
- type: variable type (NUM or CHAR);
- length: variable length (e.g. for NUM the length may be 8);
- format: variable format;

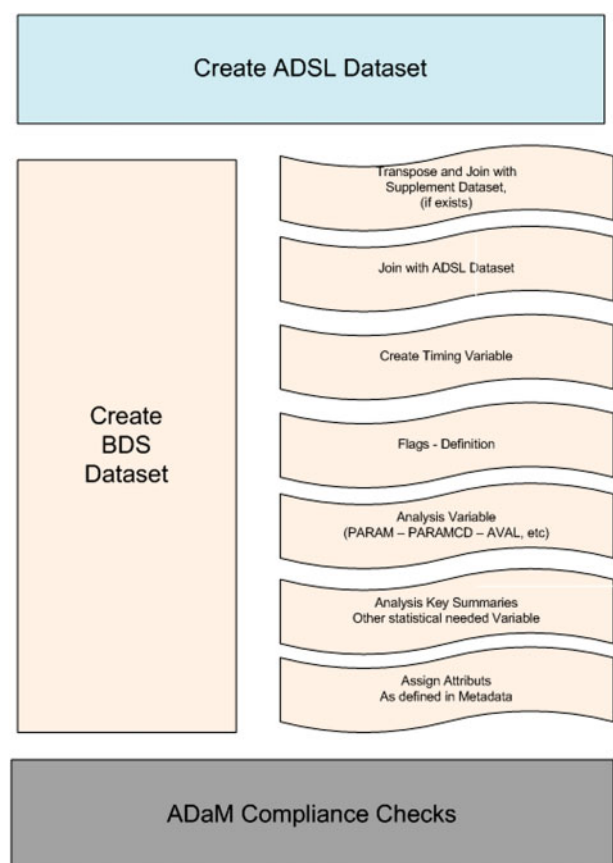


Figure 2 Create ADaM datasets

- source table: library and the table name of the source data;
- source variable: source variable (e.g. EGSTRESN from the EG domain in SDTM is the source for AVAL in the ADEG ADaM dataset);
- derivation method: describes the derivation method of the particular variable.

Where variables in the ADaM dataset have more than one source variable, we can use a special character such as # or § to delimit the list of source variables.

The metadata may also contain variables required for basic and special data checks such as primary key and 'not null' checks and others. Such data checks do not normally fall within the statistical programming remit, yet they are quick and easy to implement within the system and can help us to identify data issues in the statistical analysis outputs. Basic data checks are described later.

Examples of analysis metadata for the automation process are shown below (Figs. 3 and 4).

A library of standard macros is used for common derivation requirements. More complex or study-specific derivations which cannot be achieved through use of the standard macros are handled by ad hoc SAS plug-ins incorporated into the output scripts. These standard derivation macros, which are also listed in the Excel spreadsheet, are automatically applied to the mapping scripts during the automation process.

An example of library of standard macros for the automation process is shown below (See Fig. 5).

Three different types of macros are used for the automation process:

1. Macros with leading '\_':  
These are 'submacros' whose sole function is to derive a variable within a data step or SQL procedure, e.g. #Char2Num → %\_help\_c2n(var = ##AUTO##);
2. Macros without leading '\_':  
These macros create a new dataset from an input dataset, e.g. #AutoJoin → %auto\_join ( intab = ##INTAB##, outtab = ##OUTTAB## );
3. Pre- and post-processing macros:  
These macros carry out additional and often more complex pre- and post-processing tasks not covered by the macros above.

## Automation Process

User-friendly and easily read by SAS, Excel provides a convenient format for storing metadata to drive automated systems. Dataset structure, derivation method, and data checks are all defined within an

Sort	Meta	Name	Label	Type	Length	Format	SourceTab	SourceVar	DerivedMeth
100	Table	ADSL	Subject Level Analysis Data				SDTM.DM		
200	Var	STUDYID	Study Identifier	char	6				
300	Var	USUBJID	Unique Subject Identifier	char	15				
400	Var	SITEID	Study Site Identifier	char	15				
500	Var	AGE	Age in AGEU at RFSTDTC	num	8				
600	Var	WEIGHT	Weight (kg)	num	8		SDTM.VS	VSSTRESN	Q: WHERE vstestcd = 'Weight' and visit = 'SCREENING'
700	Var	HEIGHT	Height (cm)	num	8		SDTM.VS	VSSTRESN	Q: WHERE vstestcd = 'Height' and visit = 'SCREENING'
800	Var	TRT01SDT	Date of First Exposure in Period 01	num	8	DATE9	SDTM.EX	EXSTDTC	#FirstDate

Figure 3 ADSL metadata (excerpt)

Sort	Meta	Name	Label	Type	Length	Format	Informat	SourceTab	SourceVar	DerivedMeth
100	Table	ADEG	Analysis Dataset for Electrocardiogram					SDTM.EG		
200	Var	STUDYID	Study Identifier	char	6					
300	Var	DOMAIN	Domain Abbreviation	char	2					
400	Var	USUBJID	Unique Subject Identifier	char	15					
500	Var	EGSEQ	Sequence Number	num	8					
600	Var	EGTESTCD	ECG Test or Examination Short Name	char	8					
700	Var	EGTEST	ECG Test or Examination Name	char	40					
800	Var	EGCAT	Category for ECG	char	30					
2100	Var	EGDTC	Date/Time of ECG	char	20					
2200	Var	EGENDTC	End Date/Time of ECG	char	20					
2300	Var	EGPTNUM	Planned Time Point Number	num	8					
2400	Var	EGTPT	Planned Time Point Name	char	50					
2500	Var	AVAL	Analysis Value	num	8			SDTM.EG	EGSTRESN	
2600	Var	CHG	Change from Baseline	num	8			DERIVED		AVAL - BASE
2700	Var	BASE	Baseline Value	num	8			SDTM.EG	EGSTRESN	Q: WHERE EGBLFL = 'Y'
2800	Var	AVAL2	Analysis Value - Converted from EGSTRES	num	8			SDTM.EG	EGSTRESC	#Char2Num
2900	Var	CHG1G	Change Group	char	8			DERIVED		IF (CHG/BASE) < 10 THEN chg1g = 1; ELSE chg1g=2;

Figure 4 ADEG metadata (excerpt)

Excel 'driver sheet' from which a set of SAS macros generates the SAS scripts that, in turn, create the ADaM datasets.

General steps during the automation process:

1. import the ADaM metadata into SAS;
2. import the functional core (macro libraries) into SAS;
3. attach the functional core definitions to the ADaM metadata;
4. the 'Definition Step' (see below);
5. create mapping scripts based on the metadata.

For example, the metadata of the ADEG dataset defines four derived variables (AVAL, CHG, BASE, and CHG1G). We know from the definitions we have specified in the metadata that AVAL and BASE need to be derived before CHG and CHG1G which are dependent upon them. This example demonstrates the importance of the 'Definition Step' in the automation process.

Accordingly, the definition step determines the order of derivation, regardless of how these variables are sorted in the metadata, by looping through each observation and performing a keyword search of the variable 'DerivedMeth' in the metadata and the Processing variable in the functional core. Thus, derived variables upon which other derived variables depend are identified so that these can be derived first.

In our example, the process creates the BASE variable as the baseline value in a separate dataset which is then merged back onto the original dataset. This enables the subsequent derivation of the change from baseline variables.

The macro call for this whole process is:

```
%ADaM_Gen (metadata=ADEG
, runall=N
, addobs=ENDPOINT # AVERAGE)
```

Name	Description	LibMacro	PreProcessing	Processing	PostProcessing
MinDate		%mindate( intab = ##INTAB## , outtab = ##OUTTAB## , outdate = ##NAME## , date = ##AUTO##):			
Char2Num		%_help_c2n(var = ##AUTO##)			
AutoJoin		%auto_join( intab = ##INTAB## , outtab = ##OUTTAB##):			
HelpSort				PROC SORT DATA = ##INTAB## OUT = ##OUTTAB##; BY USUBJID ##BY##;	
FirstDate			#HelpSort	#MinDate	

Figure 5 Example of macro libraries

The macro %ADaM\_Gen generates the ADEG mapping script which produces all of the variables defined in the metadata and creates derived observations with the calculated value of DTYPE='ENDPOINT' and 'AVERAGE' (See Figs. 6 and 7).

### Validation — Data Checks

The system automatically creates ADaM datasets by producing a series of SAS scripts, one per ADaM dataset, which can then be validated against the original ADaM dataset specifications. This ensures that validation can take place in full compliance with SOPs. And because program structure is controlled, syntax and logic errors are avoided, consistency across studies and projects is achieved and validation further facilitated.

The Excel driver sheet may include basic data checks such as cross-dataset referencing to ensure completeness of the data and checks to ensure that

values are within range or not null. A data error report is generated and this can be passed to Data Management to facilitate data cleaning. This should cover only basic quality checks such as primary key and missing value checks, but simple cross checks with other datasets are also possible (See Fig. 8).

Further SAS macros to ensure that the generated ADaM Datasets are compliant with the ADaM Implementation Guide can be built into the automation macro and scheduled to run after the ADaM dataset is created.

### Conclusions

Long experience of clinical programming and sound knowledge of the SDTM and ADaM Implementation Guides are key to the successful development of a system like this. Once the automation system is created though, much less time and effort are needed to define the Excel metadata according to the SAP,

```

/*****
* Program name       : ADEG_script.sas
* (This script is automatic generated by %ADaM_Gen
* Author            : Endri
* Date created      : 29.07.2011
* Study             : (Study number)
*                   : (Study title)
* Purpose           : ADEG - Analysis Dataset for Electrocardiogram
* Template          :
* Inputs            :
* Outputs           :
* Program completed : Yes/No
* Updated by        : (Name) - (Date):
*                   : (Modification and Reason)
*****/

/* Analysis Value # Analysis Value - Converted from EGSTRESC - SDTM.EG */
DATA adeg_010_eg;
  SET sdtm.eg;
  aval = egstresn;
  aval2 = %_help_c2n(var = egstresc);
RUN;

/* Baseline Value - SDTM.EG */
DATA adeg_020_base;
  SET sdtm.eg;
  base = egstresn;
  WHERE EGBLFL = 'Y';
RUN;

/* Join all */
%auto_join ( intab = adeg_010_eg
             $ adeg_020_base # base
             , outtab = adeg_30_all);

/* Change from Baseline # Change Group - DERIVED */
DATA adeg_40_chg;
  SET adeg_30_all;
  chg = AVAL - BASE;
  IF (CHG/BASE) < 10 THEN chg1g = 1;
  ELSE chg1g = 2;
RUN;

/* Adding derived observation - DERIVED*/
%calc( intab = adeg_40_chg
      , outtab = adeg_50_calc
      , calc = ENDPOINT # AVERAGE);

/* Assign label and data check */
%assign_attrih( metadata = ADEG
               , intab = adeg_50_calc);

```

Figure 6 ADEG — SAS script (excerpt written by the automation macro)



```

/*****
* Program name      : ADSL_script.sas
* (This script is automatic generated by %ADaM_Gen
* Author           : Endri
* Date created      : 29.07.2011
* Study            : (Study number)
*                  : (Study title)
* Purpose          : ADSL - Subject Level Analysis Data
* Template         :
*
* Inputs           :
* Outputs          :
* Program completed : Yes/No
* Updated by       : (Name) - (Date):
*                  : (Modification and Reason)
*****/

/* SDTM.DM */
DATA adsl_010_dm;
  SET sdtm.dm (KEEP = studyid usubjid siteid age);
RUN;

/* Weight (kg) - SDTM.VS */
DATA adsl_020_weight;
  SET sdtm.vs;
  weight = vsstresn;
  WHERE vstestcd = 'WEIGHT' and visit = 'SCREENING';
RUN;

/* Height (cm) - SDTM.VS */
DATA adsl_030_height;
  SET sdtm.vs;
  height = vsstresn;
  WHERE vstestcd = 'HEIGHT' and visit = 'SCREENING';
RUN;

/* Date of First Exposure in Period 01 - SDTM.EX */
PROC SORT
  DATA = sdtm.ex
  OUT = adsl_040_tr01sdt;
  BY usubjid exstdtc;
RUN;

%mindate( intab = adsl_040_tr01sdt
, outtab = adsl_041_mindate
, outdate = tr01sdt
, date = exstdtc);

/* Join all */
%auto_join ( intab = adsl_010_dm
              $ adsl_020_weight # weight
              $ adsl_030_height # height
              $ adsl_041_mindate # tr01sdt
, outtab = adsl_50_all);

/* Assign label and data check */
%assign_attrik( metadata = ADSL
, intab = adsl_50_all);

```

Figure 7 ADSL — SAS script (excerpt written by the automation macro)

and manual programming is only needed for more complex tasks.

Other advantages of this system are:

- the automatically created mapping scripts are standard so that double programming is no longer required;
- documentation of the analysis datasets is already done and only needs adjustment if required;
- an ‘algorithm template’ library can be built up from complex solutions for reuse in subsequent projects;
- basic data checks for Data Management are easily built in as needed;
- the system ensures logical and consistent naming of intermediate datasets within the scripts produced, and provides proper descriptive comments throughout the code (two aspects of good programming practice which are all too easily neglected when coding by hand!);

- checks that ensure compliance with the Implementation Guide are also provided at the end of the mapping process.

So can all ADaM datasets for an entire study be defined and produced within a day? Of course that depends greatly on the amount of non-standard programming required, but setting up standard definitions, and generating and validating the resulting scripts is certainly achievable in a day once familiarization with the system is achieved.

### Acknowledgements

The authors would like to thank Keith Falconer (Principal Statistical Programmer, i3) for his input and thoughtful review of this manuscript.

Sort	Meta	Name	Label	ManualCheck
100	Table	ADEG	Analysis Dataset for Electrocardiogram	
200	Var	STUDYID	Study Identifier	
300	Var	DOMAIN	Domain Abbreviation	
400	Var	USUBJID	Unique Subject Identifier	DM.USUBJID
500	Var	EGSEQ	Sequence Number	
600	Var	EGTESTCD	ECG Test or Examination Short Name	
700	Var	EGTEST	ECG Test or Examination Name	
800	Var	EGCAT	Category for ECG	
2100	Var	EGDTC	Date/Time of ECG	
2200	Var	EGENDTC	End Date/Time of ECG	
2300	Var	EGTPTNUM	Planned Time Point Number	
2400	Var	EGTPT	Planned Time Point Name	
2500	Var	AVAL	Analysis Value	
2600	Var	CHG	Change from Baseline	
2700	Var	BASE	Baseline Value	
2800	Var	AVAL2	Analysis Value - Converted from EGSTRESC	
2900	VAR	CHG1G	Change Group	

Figure 8 ADEG metadata — manual check (excerpt)

## References

- 1 CDISC. Analysis Data Model (ADaM) Implementation Guide [document on the Internet]. Version 1.0. 17; Round Rock, TX: CDISC; 2009. Available from: <http://www.cdisc.org/adam>. [Accessed 15 March 2011].
- 2 CDISC. Analysis Data Model (ADaM) [document on the Internet]. Version 2.1. 17; Round Rock, TX: CDISC; 2009. Available from: <http://www.cdisc.org/adam>. [Accessed 15 March 2011].

Copyright of Pharmaceutical Programming is the property of Maney Publishing and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.