

#197

# Documenting Electronic Data Files and Statistical Analysis Programs

## Guidance for Industry

### *Draft Revised Guidance*

*This guidance is being distributed for comment purposes only.*

*This version of the guidance replaces the version made available December 2015. The document has been revised to update contact information, clarify existing language, remove recommendations that are no longer applicable, and provide additional details on the README file.*

Submit comments on this draft revised guidance by the date provided in the Federal Register notice announcing the availability of the draft guidance. Submit electronic comments on the guidance at <https://www.regulations.gov/>. Submit written comments to the Dockets Management Staff (HFA-305), Food and Drug Administration, 5630 Fishers Lane, Room 1061, Rockville, MD 20852. All comments should be identified with the Docket No. FDA-2009-D-0052.

For further information regarding this document, contact [Virginia Recta](#), Center for Veterinary Medicine (HFV-160), Food and Drug Administration, 7500 Standish Place, Rockville, MD 20855, 240-402-0840, email: [virginia.recta@fda.hhs.gov](mailto:virginia.recta@fda.hhs.gov).

Additional copies of this draft revised guidance document may be requested from the Policy and Regulations Staff (HFV-6), Center for Veterinary Medicine, Food and Drug Administration, 7500 Standish Place, Rockville, MD 20855, and may be viewed on the Internet at either <https://www.fda.gov/AnimalVeterinary/default.htm> or <https://www.regulations.gov/>.

**U.S. Department of Health and Human Services  
Food and Drug Administration  
Center for Veterinary Medicine  
May 2018**

**Table of Contents**

- I. INTRODUCTION..... 3**
- II. BACKGROUND ..... 3**
- III. STRUCTURE AND CONTENT OF README FILES ..... 4**
  - A. The README file ..... 4
  - B. Structure of README Files..... 5
    - 1. Electronic Data Files ..... 5
    - 2. Audit Trail Files ..... 6
    - 3. Data Analysis Programs..... 8
- IV. DATA AND ANALYSIS PROGRAMS: ADDITIONAL COMMENTS..... 10**
- APPENDIX..... 11**

## **Documenting Electronic Data Files and Statistical Analysis Programs**

### **Draft Revised Guidance for Industry**

*This draft revised guidance, when finalized, will represent the current thinking of the Food and Drug Administration (FDA or Agency) on this topic. It does not establish any rights for any person and is not binding on FDA or the public. You can use an alternative approach if it satisfies the requirements of the applicable statutes and regulations. To discuss an alternative approach, contact the FDA staff responsible for this guidance as listed on the title page.*

#### **I. INTRODUCTION**

This guidance is provided to inform sponsors of recommendations for documenting electronic data files and statistical analyses submitted to the Center for Veterinary Medicine (CVM) to support new animal drug applications. These recommendations are intended to reduce the number of revisions that may be required for CVM to effectively review data submissions. They are also intended to simplify submission preparation for sponsors by providing a suggested documentation framework, including a sample structure on how to describe and organize the information regarding the electronic data files and statistical analysis programs.

The determination of what data and analysis are needed to support a new animal drug application may vary depending on many factors and is outside the scope of this document. You should refer only to those portions of this guidance that are applicable to your particular submission.

In general, FDA's guidance documents do not establish legally enforceable responsibilities. Instead, guidances describe the Agency's current thinking on a topic and should be viewed only as recommendations, unless specific regulatory or statutory requirements are cited. The use of the word *should* in Agency guidances means that something is suggested or recommended, but not required.

#### **II. BACKGROUND**

For new animal drug applications, FDA requires full reports of investigations which have been conducted to show a drug is safe and effective for use [section 512(b)(1)(A) of the Federal Food, Drug, and Cosmetic Act (the FD&C Act)]. Additionally, section 512(n)(1)(E) of the FD&C Act requires that abbreviated applications for the approval of a new animal drug contain information to show that the new animal drug is bioequivalent to the approved new animal drug.

Submissions to CVM in support of new animal drug applications generally include a Final Study Report (FSR) for one or more studies. For each study that includes electronic data files, CVM needs information regarding documentation of the process for data generation and the statistical analysis conducted in order to review submissions and verify that there is sufficient quality and detail of evidence to support an animal drug application. An adequately documented submission

## ***Contains Nonbinding Recommendations***

### ***Draft – Not for Implementation***

should include readable electronic data files, a description of how the data are processed, and a description of the statistical analyses employed to support your conclusions.

Your documentation should clearly describe the entire process by which the data were collected, including a record of all changes to the data, starting from the electronic data files created from transcribed case report forms, or from electronic data capture systems, to the completed statistical analyses which form the basis for your study's conclusions. To understand the process by which you compiled the data and conducted the statistical analysis, CVM needs to understand the contents of each data file, the computer programs that processed all the electronic data files for analysis, and the programs that implemented the statistical analyses. The information that should be submitted to CVM, together with the FSR and electronic datasets, are described in Sections III and IV. Section III describes recommendations for how README files should be organized and completed to describe the datasets and analysis, and Section IV contains additional recommendations regarding statistical programs.

### **III. STRUCTURE AND CONTENT OF README FILES**

Data submissions to CVM typically include data in Statistical Analysis System (SAS) transport XPORT (XPT) or non-proprietary eXtensible Markup Language (XML) format, analysis programming files in XML format, and documents in Portable Document Format (PDF). A manual or instructional resource that contains important information about the electronic data files in the submission should be submitted to CVM in a README file. In this section, we describe the information that should be in the README file for CVM to evaluate the electronic data files.

#### **A. The README file**

An overview of electronic data files, documentation, and programming files included in the submission helps CVM assess the submission and the files for review. The README file contains information about the electronic data files and programming files in a submission. The README file explains how the datasets are organized and describes the programs used for dataset generation and data analysis. An effective README file should quickly orient the user to crucial information needed to understand the electronic files in a submission.

The README file is typically a PDF file with the filename README.pdf. The README file should not contain data, interpretation of the data, literature, references, notes to file, protocol-associated documents, communication records, personnel records, information not needed to interpret the data provided in the submission, or other information needed for the technical section.

A submission may contain one or more README files depending on the organization and complexity of the submission. You may choose to combine information for all studies contained in a submission into one README file, or provide separate README files for each study. README file(s) should be separate from the FSR.

*Contains Nonbinding Recommendations*  
*Draft – Not for Implementation*

**B. Structure of README Files**

The README file should include a brief introduction that includes the study number or other identifier, and study descriptor along with general orientation, background and other information relevant to analyzing and interpreting the data for each study. A description of data flow could also be included, such as audit trail processes or how the data were captured and merged to derive the analysis datasets.

The following describes the README file and a sample structure.

**1. Electronic Data Files**

a. List of Data Files

In this section, you should provide a table listing the electronic data files submitted in XML or XPT file format with brief descriptions. This table should include the file name, a brief description of contents, name of data collection form if applicable, and information on how the data were collected or any reference to location of information in the FSR that is needed to interpret the data (e.g., reference ranges or scoring definitions). See Example Table 1a.

**Example Table 1a. Listing of XML or XPT Data Files**

<b>File Name<sup>1</sup></b>	<b>File contents</b>	<b>Name of Data Collection Form (e.g., Case Report Form)</b>	<b>Comments<sup>2</sup></b>
ClinObs.xml	Daily clinical observations during hospitalization	Observation Form	Clin Obs were collected via [name of data capture system]
OwnerObs.xpt	Owner observation results for individual animals	Owner Diary	Recorded manually and transcribed.
PlateletMan.xpt	Secondary variable: Manual platelet count	Clinical Pathology Form	Reference ranges available in Appendix G of Final Study Report

<sup>1</sup> File name extension included in order to identify the file as XML or XPT.

<sup>2</sup> Information included on how data were collected (paper data collection forms or electronic data capture (EDC) system) or any reference(s) to location of information in the FSR (e.g., reference ranges or scoring definitions) that are needed to interpret the dataset.

**Contains Nonbinding Recommendations**  
*Draft – Not for Implementation*

b. Data File Contents

For each data file submitted, you should provide a table that includes the variable names, the abbreviations used in the file, variable label or description, formulas for derived variables, and additional details (e.g., description of coded values, unit of measure, formatting information), if applicable. Results from the CONTENTS procedure in SAS are not sufficient. If values were computed, derived, or transformed from other variables, the equation(s) for each variable and a table of the calculated values should be in the FSR. The calculation can be briefly described. See Example Table 1b. Standardized internal file formats for clinical and nonclinical study data are acceptable [e.g., Clinical Data Interchange Standards Consortium-Study Data Tabulation Model (SDTM) and Standard Exchange for Nonclinical Data (SEND)].<sup>1</sup>

**Example Table 1b. Contents of DataFileName.XML (number of observations=, number of variables=)**

<b>Variable Name</b>	<b>Variable Label or Description</b>	<b>Additional Details (e.g., descriptions of coded values, units of measure, formatting information)</b>
Animal	Animal Identification	
Dot	Day of treatment	Study day formatted as mm/dd/yyyy
Trt	Treatment	T01=Placebo; T02=Drug
Dscore	Daily Depression Score	0=normal; 1=slight; etc.
OverallScore	Overall Depression Score	Sum of Daily Depression scores (Day 1 to 5)

**2. Audit Trail Files**

This section only applies to studies that include audit trail files as part of the controls that ensure the authenticity and integrity of records in electronic data capture (EDC) systems.<sup>2</sup> The audit trail is a portion of the raw data<sup>3</sup> that includes the original recorded data point and any changes to the data point, the identification of individuals that entered or changed data in the EDC system, the date and time the data point was entered or changed, the reason for each change, and the date when the database was

<sup>1</sup> See Study Data Standards: <https://www.fda.gov/forindustry/datastandards/studydatastandards/default.htm>.

<sup>2</sup> Electronic Data Capture (EDC) or Electronically Captured Data (ECD) refers to data that was electronically captured at first observation (e.g., web-based software or analytical instrument).

<sup>3</sup> As defined in 21 CFR 58.3(k) and in FDA Guidance for Industry #85 (VICH GL 9), “[Good Clinical Practice](https://www.fda.gov/downloads/AnimalVeterinary/GuidanceComplianceEnforcement/GuidanceforIndustry/UCM052417.pdf)” (<https://www.fda.gov/downloads/AnimalVeterinary/GuidanceComplianceEnforcement/GuidanceforIndustry/UCM052417.pdf>).

## ***Contains Nonbinding Recommendations***

### ***Draft – Not for Implementation***

locked.<sup>4</sup> If you are submitting audit trails on specified variables, the electronic audit trails should be submitted in XPT or non-proprietary XML file format.

CVM's preference is to receive one copy of the EDC study database that includes the final variable observations (e.g., the electronic data files described in III.B.1 above) and associated audit trail information, described in this section, together in one electronic data file or set of data files. If you submit these combined data sets, you should also submit the programs needed to create data files that contain the final observations suitable for review and statistical analysis.

#### a. Audit Trail File Listing

If you are submitting the electronic audit trail separately from the data files for review and statistical analysis, you should provide a table listing the audit trail files submitted. This table should include the file name, the description of the file including EDC system name, and any information necessary for review. See Example Table 2a.

#### **Example Table 2a. Listing of Audit Trail Files.**

<b>File Name</b>	<b>Description</b>	<b>Comments</b>
Site_A_Audit.xml	Audit trail of Site A	[name of EDC system]
Site_B_Audit.xml	Audit trail of Site B	[name of EDC system]

#### b. Audit Trail File Contents

For each audit trail file submitted, you should provide a table that includes the variable names, variable label or description (e.g., description of coded values, unit of measure, formatting information), and other information necessary for review. See Example Table 2b. Each audit trail file submitted should include the original and updated values of each data point, operator identification and date and time stamps for each data entry and any change, and reason for each change. Additional columns may be added as needed.

---

<sup>4</sup> See Guidance for Industry “[Computerized Systems Used in Clinical Investigations](https://www.fda.gov/downloads/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/UCM070266.pdf)” (<https://www.fda.gov/downloads/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/UCM070266.pdf>) or Guidance for Industry “[Part 11, Electronic Records; Electronic Signatures - Scope and Application](https://www.fda.gov/downloads/RegulatoryInformation/Guidances/ucm125125.pdf)” (<https://www.fda.gov/downloads/RegulatoryInformation/Guidances/ucm125125.pdf>).

*Contains Nonbinding Recommendations*

*Draft – Not for Implementation*

**Example Table 2b. Contents of AuditTrailFile.XML (number observations=, number of variables=)**

<b>Variable Name</b>	<b>Description</b>
animalid	Animal identification number
formname	Name of the data collection form
entryfield	Name of entry field from data collection form
entrytype	Type of entry field (e.g, numeric, text, radio button or checkbox)
entry	Data entry [include explanation if applicable (e.g., “on” = radio button was selected, “off” = radio button not selected; if blank, no data entered)]
operatorid	Individual entering data
entrydate	Date/time stamp when the data was entered
reason	Reason for change, if applicable
lockdate	Date and time when the database was locked

**3. Data Analysis Programs**

a. Program File Listing

In this section, you should provide a table listing the programs used to perform randomization, process the data, generate summaries, and perform the statistical analysis. This table should include the file name, the purpose of the program, the electronic data files accessed and generated by each program, and a list of any results (tables/graphs) generated, if applicable. See Example Table 3.

**Example Table 3. Listing of Program Files.**

<b>Program File Name</b>	<b>File Description</b>	<b>Data File(s) Used in Program (Input)</b>	<b>Data File(s) Created (Output)</b>	<b>Results or Tables/ Graphs Created</b>	<b>Comments</b>
setup.xml	Set up needed libraries and data formats. Creates all temporary data files used in the analysis	Effect_data.xpt AE_all.xpt Clinobs_all.xpt	ClinObs.sas7bdat Effect.sas7bdat AE.sas7bdat	None	Run this program first before running any analysis programs

***Contains Nonbinding Recommendations***  
*Draft – Not for Implementation*

<b>Program File Name</b>	<b>File Description</b>	<b>Data File(s) Used in Program (Input)</b>	<b>Data File(s) Created (Output)</b>	<b>Results or Tables/ Graphs Created</b>	<b>Comments</b>
s_clinobs.xml	Summarize clinical observations	ClinObs.sas	None	Summary of clinical observations  Profile plots	Table X provided in Final Study Report
effect_prog.xml	Summarize effectiveness variables and determine success/failure	Effect.sas7bdat	Succes.sas7bdat	None	Details of success/failure evaluation. Table Y provided in Final Study Report.
prim_anlys.xml	Conduct primary analysis	Succes.sas7bdat	None	Primary analysis results	

b. Overview of Data and Analysis Flow

You should provide a general overview of the data and analysis process used in the study and submission. For example, you should describe whether data sets were directly downloaded from data management systems and whether any processing (merging, validation, editing) were performed prior to analysis.

Because databases may not store or directly export data in XPT or non-proprietary XML formats, conversion from another file format, e.g., XLSX, CSV, or SAS7bdat, may be needed. If conversion is necessary, CVM expects that electronic data files will be converted into XPT or non-proprietary XML format without modification (e.g., no column changes or calculations) using validated software. The Appendix has examples of SAS programs that convert data files to acceptable formats. The program used to make the conversion should be provided to CVM in XML format. Additionally, CVM encourages you to evaluate/analyze the data using the same file that is submitted to CVM. For example, if the study data are submitted to CVM in file “ALLDATA.XML,” the programs submitted should use this file for analysis.

Submitted datasets should be unmodified after export from the database. If a dataset was modified during analysis (e.g., variable or animal excluded or coded information reformatted as text), a description of the changes and the computer program used to create the changes should be included, but the modified dataset files should not be submitted. The logic or coding that was used in creating the modified datasets should be documented in the computer program and clearly detailed in the submission so that CVM may recreate and verify the modifications.

*Contains Nonbinding Recommendations*  
*Draft – Not for Implementation*

c. Instructions for Running Programs

You should describe the sequence of program calls needed for CVM to run your programs. Starting with the first program to be run, you should describe calls to other programs, custom functions, and macros, if any.

For each program, you should document all directories and files referenced to access or store data, including directory and file names, locations, and aliases if used. Describe programs defining custom styles or formats or, if such styles or formats are predefined, you should provide instructions for their installation. If the programs were designed to call other programs or access data in specific folders or directory structures, you should describe this structure so that CVM can verify your analysis process.

d. Randomization Programs

You should list and submit any programs used to generate random assignments, such as allocation of experimental units to treatment groups or determining order of necropsy. The randomization programs should include randomization seed(s), if used.

The details on the randomization process, the actual programs used, and resulting allocation tables should be included in the FSR.

#### **IV. DATA AND ANALYSIS PROGRAMS: ADDITIONAL COMMENTS**

It is acceptable to submit a separate document (e.g., Statistical Report) as an appendix in the FSR that provides details on the statistical analysis as well as additional analysis results and summaries. If the Statistical Report or other document includes the information regarding program files and program execution as described in III.B.3 above, then the information does not have to be repeated in the README file. Instead, the README file can refer to the appropriate sections of other document(s) in the submission.

CVM's verification of the analysis process will be facilitated by submission of well-documented analysis programs. You should describe the general purpose of each program in the beginning of the statistical program. Within the program, you should include sufficient comments to explain complex sections, for example, if a program will call certain macros or subroutines.

You should not submit electronic files generated by the statistical analysis, such as tables of descriptive summaries or least squares means, in data format (e.g., XML or XPT). Documents containing these summaries may be included or appended to the FSR as appropriate. For example, a listing of subject success/failure outcomes should be included in the FSR but the secondary electronic data file derived from the original raw data should not be submitted. Additionally, you should not submit outputs from running the analysis; for example, the outputs and log files produced when running SAS.

## **APPENDIX**

Examples of SAS codes that convert data files to acceptable formats

1. SAS code to generate non-proprietary XML files

```
libname in 'file location';  
libname out xml 'file location\filename1.xml';  
data out.dataset1; set in.filename; run;
```

2. SAS code to generate XPT files

```
libname in 'file location';  
libname out xport 'file location\filename1.xpt';  
data out.filename1; set in.filename; run;
```