

Guideline on the Submission of Clinical Trial Data

(Draft for Public Review)



Center for Drug Evaluation, NMPA
May, 2020

Table of Contents

1. Background and Purposes	3
2. Submission Components of Clinical Trial Data	4
2.1 Study database	4
2.2 Analysis database	5
2.3 Data Definition File	6
2.4 Annotated CRF	7
2.5 Programming Code	7
3. Submission Document Format and Conventions	8
3.1 Portable document format	8
3.2 Extensible mark-up language format	8
3.3 Plain text format	8
3.4 Data transport file format	8
3.5 Dataset split	9
3.6 Dataset name, variable name and length	9
3.7 Dataset labels and variable labels	10
4. Other Considerations	10
4.1 Traceability of trial data	10
4.2 Data files under eCTD	11
4.3 Foreign language database	11
4.4 Communication with regulatory agency	11
References	11
Appendix 1: Commonly Used Raw Datasets	13
Appendix 2: STF	14
Appendix 3: Folder structure	15
Appendix 4: Glossary	16

Guideline on the Submission of Clinical Trial Data

1. Background and Purposes

Clinical trial data is one of the important materials submitted by sponsor to regulatory agencies. It is a valuable resource for both regulatory agencies and sponsor. Standardized collection, organization, analysis, and presentation of clinical trial data play an important role in improving the efficiency and quality of clinical researches and development, shortening review timelines. Also, it is beneficial to the management of entire life cycle of drug development, and further promoting information exchange and sharing between drug development units and regulatory agencies.

If the clinical trial data submitted by sponsor does not follow certain specifications, it will take significant resources for reviewers to get familiarized with and understand the data structure and content. In some cases, sponsor or regulatory agency may need to conduct pooled analyses using multi-sources of clinical trial data. Non-standardized data will make this task almost impossible.

Submission of clinical trial data is usually a package that includes database and its supportive documents, like data definition file, data reviewer's guide, data derivation and analysis programs, and annotated case report forms (aCRFs). This document provides specific requirements for the content and format of clinical trial data submission package, and aims to guide sponsors in the submission of clinical trial data and related materials, and to help professionals such as data managers and statistical programmers to better conduct related tasks for clinical trials.

This guideline is formulated based on the data submission requirements of international regulatory agencies and the current situation in China. Sponsor should prepare the package based on the requirements in this guidance document. Sponsor is encouraged to submit clinical trial data and the associated materials according to Clinical Data Interchange

32 Standards Consortium (CDISC) standards. With the development and improvement of the
33 understanding and practice of clinical trial data standards, this guidance document will be
34 revised and improved as appropriate.

35

36 **2. Submission Components of Clinical Trial Data**

37

38 2.1 Study database

39

40 Study database should generally contain source data collected directly from case report
41 forms (CRFs) and external sources as well. It may also contain few derived variables such
42 as serial numbers, study days, etc. However, missing data in raw database should not be
43 imputed. To meet the requirements for regulatory submission, collected data may require
44 necessary standardization or coding (e.g., adjusting the name/label of the dataset in the
45 database and name/label of the variable in the dataset, encoding variable values with
46 standard dictionary (e.g. Medical Dictionary for Regulatory Activities (MedDRA)) where
47 applicable.

48

49 Study database typically contains multiple raw datasets, which should be organized and
50 named according to its contents. Study dataset is usually named as a two letter code, for
51 example, demographics dataset (dm), adverse event dataset (ae), laboratory test dataset (lb),
52 etc. Refer to Appendix 1 for details of nomenclature of raw datasets commonly used in
53 clinical trial data submission.

54

55 Datasets that contain observed results of subjects must have study identifier, subject unique
56 identifier, and other identifiers to uniquely identify an observation must be included in a
57 dataset (e.g., dm, ae, lb, etc. in Appendix 1). Subject identifier (SUBJID) must be included
58 in dm dataset. Commonly used identifiers are exemplified as follows:

59 Study identifier: The variable name is STUDYID, character type, unique identifier of the
60 study, usually also regarded as study number.

61

62 Subject unique identifier: The variable name is USUBJID, character type. Each subject
63 should be assigned a unique identifier throughout a submission (may include multiple
64 clinical studies). In all datasets (including the raw and analysis datasets), the same subject
65 should have exactly the same unique identifier. When a subject participates in multiple
66 studies, the USUBJID should be consistent across these studies. Following this rule is
67 particularly important for merging datasets from different studies on the same subject (e.g.,
68 randomized controlled study and its extension study).

69

70 Subject identifier: The variable name is SUBJID, character type. SUBJID is the identifier
71 of a subject enrolled in a trial. If one subject is screened multiple times in a trial, the
72 SUBJID should be different each time.

73

74 Time variables such as visit name (VISIT, character type) and visit number (VISITNUM,
75 numerical type) should be included in applicable datasets. VISITNUM should be assigned
76 values in ascending chronological order.

77

78 2.2 Analysis database

79

80 Analysis database is a database derived for statistical analyses and is used to produce and
81 support statistical analysis results in clinical study report. Analysis database can contain
82 raw data, and also derived variables according to specified rules, such as imputation for
83 missing values.

84

85 Analysis database typically includes multiple analysis datasets. Derived and collected data
86 (from raw datasets or other analysis datasets) may be combined into a single dataset when
87 building an analysis dataset. When creating analysis dataset, the following principles
88 should be followed: 1) Analysis dataset must be built up to clearly support the statistical
89 analyses planned for the clinical study. 2) Analysis data set must be traceable; and the
90 specified rules for derived variables should be detailed in the corresponding data definition
91 file. 3) The structure and contents of analysis dataset should facilitate statistical analysis
92 with limited programming efforts, namely analysis ready.

93

94 Analysis database should contain all variables required for the planned/intended analysis,
95 including derived variables; and all derived variables should be able to be generated from
96 the study database. Analysis datasets are usually named in the form of "adxxxxxx" and the
97 name should be also consistent with the corresponding raw dataset's name, such as adcm,
98 adae, adlb, etc.

99

100 The subject level analysis dataset is mandatory (named as adsl) for a submission data
101 package. In this dataset, each subject should have one record that includes, but not limited
102 to, demographics, disease factors, treatment groups, other prognostic factors that may
103 affect treatment response, dates of important events, and population flags.

104

105 For some endpoints (e.g., scale scores), a series of derivation processes are needed to get
106 it ready for the final statistical analyses. The intermediate variables/datasets derived to
107 facilitate the creation of the final analysis dataset should also be included in the analysis
108 database for submission, if necessary.

109

110 2.3 Data Definition File

111

112 Raw and analysis databases submitted must have appropriate data definition file. Data
113 definition file is used to describe the submitted data, and should at least contain the name,
114 label and basic structure of each dataset in the submitted database, and the name, label and
115 type of each variable and derivation process of each derived variable in each dataset.

116 Data definition file is one of the most important documents for regulatory agencies to
117 accurately understand the submission data. Sponsor should ensure that the code list and
118 sources of each variable are clearly defined and easily searchable. If external dictionaries
119 are used, sponsor needs to specify the dictionary and its version in data definition file. Good
120 traceability between data (e.g., between raw data and CRF, analysis data and raw data)
121 needs to be documented in the file to facilitate regulatory review. Sponsor need to provide
122 details in data definition file, particularly with regard to derived variables. Program codes
123 may need to be provided, if necessary, to assist with review.

124

125 Data reviewer's guide is a supplement to data definition file for raw/analysis database,
126 which will help reviewers better understand and use submitted data, so it should be
127 submitted if necessary. Data reviewer's guide provides information in addition to what are
128 presented in data definition file, including but not limited to, instructions on the use of the
129 submitted data, relationships between the study report and the data, certain key information
130 of study documents (e.g. trial protocol, statistical analysis plan, clinical study report), and
131 description/explanation of other special scenarios. Data reviewer's guides are not intended
132 to replace data definition file, but to help reviewers more accurately and efficiently
133 understand and use the submitted database, relevant terminologies, and data definition file.

134

135 Data definition file is generally in extensible mark-up language format (XML) or portable
136 document format (PDF) format. Data reviewer's guide should be submitted in PDF.

137

138 2.4 Annotated CRF

139

140 Annotated CRF is blank CRF with annotations that illustrates the mapping relationship
141 between data units (i.e. field) of collected subject data (electronic or paper) and
142 variables/variable values in submitted study database. Annotated CRF should be
143 submitted in PDF.

144

145 In practice, some data fields may be collected on the CRF but not included in submitted
146 datasets. These data fields should be clearly marked as "NOT SUBMITTED" on the aCRF
147 and reason(s) for not submitting these data should be clarified in data reviewer's guide
148 accordingly.

149

150 2.5 Programming Code

151

152 Sponsors should submit programming codes, which include, but not limited to, the
153 derivation process of analysis datasets, generation process of analysis results for the
154 primary and secondary efficacy endpoints, etc. Programming codes submitted in

155 submission package.should be readable (with comments), understandable, executable, and
156 do not include external program calls, which in particular avoid using large macro
157 programs. Programming codes in submission packages are generally in TXT format.

158

159 **3. Submission Document Format and Conventions**

160

161 3.1 Portable document format

162

163 Portable Document Format is an open document format that is independent of application
164 software, hardware, and operating system. Any other documents in submission package
165 that follow the requirements of the International Counsel for Harmonization (ICH)
166 Electronic Common Technical Document (eCTD) format can be in PDF format. It is
167 recommended that PDF version1.4 and above to be used in submission. All PDF files
168 should use .pdf as the file extension.

169

170 3.2 Extensible mark-up language format

171

172 Extensible Mark-up Language is a type of data exchange languages, which is defined by
173 the World Wide Web Consortium (W3C). It can be opened, edited and created by any text
174 editor, and used to transfer and store data. Files in XML format can conveniently exchange
175 information between different systems. All XML files are required to use .xml as the file
176 extension.

177

178 3.3 Plain text format

179

180 Plain Text Format document (TXT) has characteristics such as simple format, small file
181 size, simple and convenient for storage. It is also a common file format supported by
182 computers and many mobile devices. All TXT files should use .txt as the file extension.

183

184 3.4 Data transport file format

185

186 Datasets in submission package are usually in transport file format (XPT). One XPT file
187 corresponds to one dataset. XPT file name needs to be consistent with the corresponding
188 dataset name. XPT files should use .xpt as the file extension, for example, ae.xpt for
189 Adverse Event (AE), cm.xpt for Concomitant Medication (CM). SAS Transport File
190 Format version 5 (referred to as XPT V5) or above is recommended as the data submission
191 format. Sponsor should ensure that submitted datasets are free from illegible contents in
192 different operating environments.

193

194 3.5 Dataset split

195

196 When a dataset in database needs to be split because the file size does not meet submission
197 requirements, detailed rules of splitting and detailed steps of merging it back should be
198 specified in data reviewer's guide to ensure that reviewer can generate the dataset same as
199 what is prior to splitting.

200

201 3.6 Dataset name, variable name and length

202

203 Specific requirements about the name and length of dataset and variable are as follows:

204 Dataset name can only contain lowercase letters and numbers and must start with a
205 lowercase letter. The maximum length of a dataset name is 8 bytes.

206

207 The variable name can only contain upper case characters and numbers, and must start with
208 a letter. The maximum length of a variable name is 8 bytes.

209

210 The length of each character variable should be set to the maximum actual value length of
211 the variable across all datasets of the same study, to effectively control the size of the file.

212 Variable length should be set not to exceed 200 bytes; variable splitting may be needed.

213 When splitting, bytes cannot be truncated, and efforts should be made to maintain the
214 integrity of statement in each splitting variable.

215

216 3.7 Dataset labels and variable labels

217

218 For ease of review, dataset labels and variable labels should be in Chinese and should not
219 exceed 40 bytes in length. If necessary, labels can contain English letters, underlines, or
220 numbers, but cannot start with numbers. In addition, labels cannot include the following
221 cases:

222

- 223 • Unpaired half-width/full-width single/double quotation marks
- 224 • Unpaired half-width or full-width brackets
- 225 • Special characters

226

227 **4. Other Considerations**

228

229 4.1 Traceability of trial data

230

231 An important part of regulatory review is accurate understanding of the source of data, that
232 is, the traceability of data. Traceability enables reviewers to understand the relationship
233 between statistical analysis results (table, listing and figures in study report), analysis data,
234 and raw data.

235

236 The traceability of data ensures that reviewers are able to accurately:

237

- 238 • understand the construction of analysis datasets
- 239 • identify records used for derived variables and the corresponding algorithms
- 240 • understand the algorithm/model of corresponding statistical results
- 241 • establish technique used to link raw data to corresponding table(s)

242

243 When submitting study database, sponsor should ensure that regulatory reviewers can use
244 the study database to derive the analysis database that is consistent with what the sponsor
245 submitted, and that analysis database can directly reproduce statistical analysis results that

246 are also consistent with what the sponsor submitted. Traceability can be supplemented by
247 providing a detailed data flowchart from the collection to the submission.

248

249 4.2 Data files under eCTD

250

251 When it comes to registration using eCTD, all documents, trial data and associated
252 supportive documents should be organized according to the specified folder structure. All
253 submitted files should be in the correct folder and tagged using the appropriate Study
254 Tagging File (STF). Refer to Appendix 2 and Appendix 3 for more information regarding
255 STF and folder structure.

256

257 4.3 Foreign language database

258

259 When it comes to registration using foreign language database, dataset label, variable label,
260 adverse events terms, generic name of concomitant medications, medical history, and
261 names of the clinical endpoints in normalized datasets (corresponding to variables in the
262 horizontal structured datasets) should be in Chinese. CRFs, aCRFs, data definition files,
263 and data reviewer's guides should also be submitted in Chinese. Chinese translations in
264 database should be consistent with all other documents in submission package.

265

266 4.4 Communication with regulatory agency

267

268 Based on specific characteristics and complexity of clinical trial data, sponsor may, if
269 necessary, communicate with regulatory authority at Pre-NDA meetings regarding the
270 clinical trial database and relevant materials to facilitate timely and accurate understanding
271 of the clinical trial data submitted by the sponsor.

272

273 **References**

274 1. CFDA: Technical Guide for Data Management in Clinical Trials, July 2016

275 2. FDA: Study Data Technical Conformance Guide, Oct 2019

- 276 3. PMDA: Revision of Technical Conformance Guide on Electronic Study Data
277 Submissions, Jan 2019
- 278 4. CDISC: Study Data Tabulation Model Implementation Guide, Nov 2018
- 279 5. CDISC: Analysis Data Model Implementation Guide, Oct 2019



280

Appendix 1: Commonly Used Raw Datasets

281

Table 1 Common raw data sets and nomenclature

Datasets	Naming	Submission Requirements
Demography	dm	Must be submitted
Medical History	mh	If applicable
Adverse Events	ae	If applicable
Prior and Concomitant Medications	cm	If applicable
Exposure	ex	If applicable
Subject Disposition	ds	If applicable
Questionnaire	qs	If applicable
Protocol Violation	dv	If applicable
Laboratory Tests	lb	If applicable
ECG	eg	If applicable
Vital Signs	vs	If applicable
Clinical Events	ce	If applicable
Physical Examination	pe	If applicable
Disease Response	rs	If applicable

282



APR
三方协调委员会

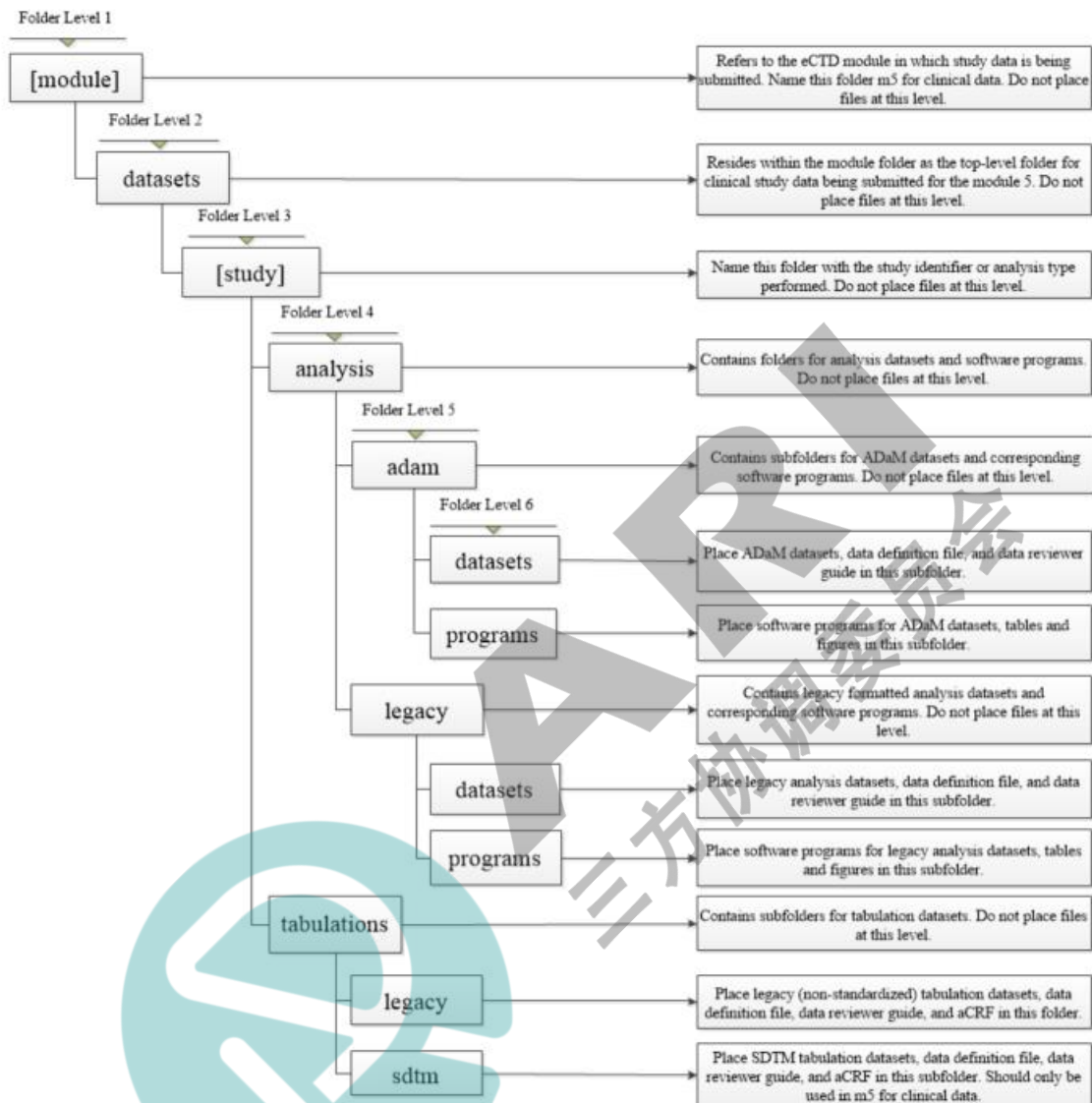
Appendix 2: STF

Name attribute values for the file-tag element	Description
data-tabulation-dataset-legacy	Study database (non-CDISC standard)
data-tabulation-dataset-sdtm	Study database (CDISC standard)
data-tabulation-data-definition	Study database data define file and data reviewer's guide
analysis-dataset-adam	Analysis database (CDISC standard)
analysis-dataset-legacy	Analysis database (non-CDISC standard)
analysis-data-definition	Analysis database data define file and data reviewer's guide
annotated-crf	Annotated CRF
analysis-program	data derivation and analysis programs



APRI
三方协调委员会

285 **Appendix 3: Folder structure**



286

287 **Appendix 4: Glossary**

288 Code List:

289 Code list for a variable is a list of allowable values that this variable may have. It includes
290 standard codes, industry commonly used codes, and sponsor custom-defined codes.

291

292 Case Report Form (CRF) :

293 A printed, optical, or electronic document designed to record all of the protocol required
294 information to be reported to the sponsor on each trial subject.

295

296 Electronic Common Technical Document (eCTD) :

297 Electronic registration documents submitted for drug registration and review. Organize,
298 transmit, and present the CTD-compliant drug submissions electronically in extensible
299 mark-up language format.

300

301 Data Definition File:

302 Data definition file is used to describe the submitted data, and should at least contain the
303 name, label and basic structure of each dataset in the submitted database, and the name,
304 label and type of each variable and derivation process of each derived variable in each
305 dataset.

306

307 Data Reviewer's Guide:

308 Data reviewer's guide is a supplement to data definition file. It includes, but not limited to,
309 instructions on the use of the submitted data, relationships between the study report and
310 the data, certain key information of study documents, and description/explanation of other
311 special scenarios.

312

313 Annotated Case Report Form (aCRF):

314 Annotated CRF is blank CRF with annotations that illustrates the mapping relationship
315 between data units (i.e. field) of collected subject data (electronic or paper) and
316 variables/variable values in submitted study database.