
Guidance for Industry

Providing Regulatory Submissions in Electronic Format — Standardized Study Data

DRAFT GUIDANCE

This guidance document is being distributed for comment purposes only.

Comments and suggestions regarding this draft document should be submitted within 60 days of publication in the *Federal Register* of the notice announcing the availability of the draft guidance. Submit comments to the Division of Dockets Management (HFA-305), Food and Drug Administration, 5630 Fishers Lane, rm. 1061, Rockville, MD 20852. All comments should be identified with the docket number listed in the notice of availability that publishes in the *Federal Register*.

For questions regarding this draft document contact (CDER) Kieu Pham at 301-796-1616, (CBER) Office of Communication, Outreach and Development (OCOD) at 301-827-1800 or 1-800-835-4709, or (CDRH) Terrie Reed at 301-796-6130.

**U.S. Department of Health and Human Services
Food and Drug Administration
Center for Drug Evaluation and Research (CDER)
Center for Biologics Evaluation and Research (CBER)
Center for Devices and Radiological Health (CDRH)**

**February 2012
Electronic Submissions**

Guidance for Industry

Providing Regulatory Submissions in Electronic Format — Standardized Study Data

Additional copies are available from:

*Office of Communications
Division of Drug Information, WO51, Room 2201
Center for Drug Evaluation and Research
Food and Drug Administration
10903 New Hampshire Ave., Silver Spring, MD 20993-0002
Phone: 301-796-3400; Fax: 301-847-8714
druginfo@fda.hhs.gov*

<http://www.fda.gov/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/default.htm>

or

*Office of Communication, Outreach and Development, HFM-40
Center for Biologics Evaluation and Research
Food and Drug Administration
5515 Security Lane, Rockville, MD 20852*

<http://www.fda.gov/BiologicsBloodVaccines/GuidanceComplianceRegulatoryInformation/Guidances/default.htm>

*(Tel) 800-835-4709 or 301-827-1800
http://www.fda.gov/cber/guidelines.htm.*

or

*Office of Communication, Education, and Radiological Programs
Division of Small Manufacturers Assistance, WO66-4613
Center for Devices and Radiological Health
Food and Drug Administration
10903 New Hampshire Avenue, Silver Spring, MD 20993-0002*

<http://www.fda.gov/MedicalDevices/DeviceRegulationandGuidance/GuidanceDocuments.htm>

Email: dsmica@cdrh.fda.gov

Fax: 301.847.8149

(Tel) Manufacturers Assistance: 800.638.2041 or 301.796.7100

(Tel) International Staff Phone: 301.827.3993

**U.S. Department of Health and Human Services
Food and Drug Administration
Center for Drug Evaluation and Research (CDER)
Center for Biologics Evaluation and Research (CBER)
Center for Devices and Radiological Health (CDRH)**

**February 2012
Electronic Submissions**

TABLE OF CONTENTS

I.	INTRODUCTION.....	1
II.	BACKGROUND	2
III.	GENERAL SUBMISSION CONSIDERATIONS	3
A.	Planning and Providing Standardized Study Data.....	4
B.	Controlled Terminologies.....	5
C.	Standardization of Previously Collected Nonstandard Data	7
D.	Data Validation	8
E.	Exceptions to Standardized Study Data Submissions	9
F.	Meetings with FDA	9
IV.	ONLINE TECHNICAL RESOURCES	10
V.	ADDITIONAL SUPPORT.....	10
	APPENDIX – DATA STANDARDS AND INTEROPERABLE DATA EXCHANGE	11
A.	Interoperability	12
B.	Types of Data Standards	14

Technical specifications and additional online resources associated with this guidance are provided separately and will be updated periodically. To ensure that you have the most recent versions, check the FDA Study Data Standards Resources Web page at <http://www.fda.gov/ForIndustry/DataStandards/StudyDataStandards/default.htm>.

Contains Nonbinding Recommendations

Draft – Not for Implementation

38
39 **These resources are available on the Study Data Standards Resources Web page at**
40 **<http://www.fda.gov/ForIndustry/DataStandards/StudyDataStandards/default.htm>.**
41

42 FDA’s guidance documents, including this guidance, do not establish legally enforceable
43 responsibilities. Instead, guidances describe the Agency’s current thinking on a topic and should
44 be viewed only as recommendations, unless specific regulatory or statutory requirements are
45 cited. The use of the word *should* in Agency guidances means that something is suggested or
46 recommended, but not required.
47

48 49 **II. BACKGROUND**

50
51 FDA routinely receives in regulatory submissions the results of scientific studies such as clinical
52 trials and animal studies. The clinical trials are generally human investigations that are intended
53 to study the pharmacokinetics, pharmacodynamics, bioequivalence, safety, and/or effectiveness
54 of an investigational product. For devices evaluated in a 510(k), the studies may compare the
55 products to a *predicate* device. Applicants typically submit study reports, which describe the
56 study protocol, the data collected, the analyses performed, the results of those analyses, and the
57 conclusions of the study. Also accompanying the study reports are the case report forms (CRFs),
58 and the study data as case report tabulations (CRTs) and analysis datasets. CRFs are the forms
59 used by the clinical investigator to document the collected data. The CRTs are aggregate (i.e.,
60 data from multiple subjects grouped together) listings of all the data collected on the case report
61 forms. CRFs and CRTs allow the Agency to perform an independent analysis of the study data.
62 FDA routinely performs its own independent analyses of study data to assess the effectiveness
63 and safety of investigational products. Analysis datasets are a subset of all data collected in the
64 study and are the critical dataset to support the primary study analyses contained in the study
65 report.
66

67 The animal studies are generally nonclinical pharmacology, toxicology, or biocompatibility
68 studies intended to investigate the actions of the investigational product in relation to its
69 proposed therapeutic indication and intended clinical use.⁴ Similar to clinical trials, applicants
70 submit study reports that include pertinent tabulations of the collected and analyzed data. FDA
71 reviewers use these data for safety and predictive assessments, which can include reanalysis of
72 the submitted data.
73

74 Existing Federal regulations allow the voluntary submission of electronic records, including
75 study data, to the Agency in lieu of paper records (see 21 CFR part 11). For many years, FDA
76 Centers have requested that clinical study data be submitted electronically because paper CRTs
77 are widely recognized as being highly inefficient to support analysis and review. The data in
78 paper CRTs are not machine-readable and therefore cannot be easily analyzed using modern
79 analytic software. Although submission of clinical study data in electronic format has become
80 relatively routine at some Centers, these data are often not standardized. An important goal of

⁴ In some circumstances, animal studies serve as efficacy studies (see 21 CFR 314.600 and 21 CFR 601.90).

Contains Nonbinding Recommendations

Draft – Not for Implementation

81 this guidance is to increase the quantity of standardized electronic study data submitted to FDA
82 for both clinical and nonclinical studies.

83
84 Standardizing study data makes the data more useful. Data that are standardized are easier to
85 understand, analyze, review, and synthesize in an integrated manner in a single study or multiple
86 studies, thereby enabling more effective regulatory decisions.⁵ As a result, FDA has publicly
87 announced that it intends to propose a new Federal regulation that would require the submission
88 of standardized electronic study data.⁶ This guidance will support the new regulation once the
89 rule is finalized.

90
91 Data that are not standardized diminish the Agency’s ability to review the data efficiently,
92 resulting in manual, labor-intensive processes and inherent inefficiencies in the review. They
93 also limit the ability to automate some routine analyses. A simple example is an analysis of the
94 distribution of the gender of subjects in a study. Gender can be represented in a study as male or
95 female, M or F, or 0 or 1 (or 1 or 2; the possible representations are endless). If sponsors always
96 expressed the gender of a subject in a study as M or F, then it is fairly trivial to create a computer
97 program that can automatically count all occurrences of M and F and automatically generate a
98 standard gender distribution report. If the next study uses 0 and 1 to represent gender, then the
99 program is no longer useful. Considering that a typical study has many hundreds of data
100 elements,⁷ many with almost limitless ways they can be represented, the challenges in analyzing
101 nonstandard data are quite evident.

102
103 As sponsors and applicants move toward standardized electronic study data submissions, there is
104 a need to understand FDA’s expectations for such data submissions. This guidance provides
105 FDA’s current thinking on the submission of standardized electronic study data.⁸

106
107 This guidance does not describe *what* data should be submitted. The content of a study data
108 submission is largely determined by science and regulation, and often involves prior discussion
109 with the appropriate review division to ensure that the content of the submission is adequate for
110 its intended purpose. *What* to submit is outside the scope of this guidance; but having once
111 decided what to submit, sponsors can consult this guidance to understand *how* to submit the data
112 electronically, using data standards supported by the Agency.

113
114 To better understand why the Agency is now emphasizing the submission of standardized data
115 for all studies, please refer to the Appendix.

116
117

118 **III. GENERAL SUBMISSION CONSIDERATIONS**

119

⁵ Standardized data also facilitate data exchange and sharing (e.g., between a contract research organization and a sponsor).

⁶ See <http://www.reginfo.gov/public/do/eAgendaViewRule?ruleID=284747>.

⁷ By *data element*, we mean the smallest (or *atomic*) piece of information that is useful for analysis. Please see the Appendix for additional information.

⁸ See section III.E for exceptions to standardized study data submissions.

Contains Nonbinding Recommendations

Draft – Not for Implementation

120 **Applicants and sponsors should refer to the FDA Study Data Standards Resources Web**
121 **page⁹ to obtain an up-to-date inventory of the various data standards supported¹⁰ by**
122 **CDER, CBER, and CDRH for standardized study data submissions.**

123
124 To fully realize the advantages and benefits of data standardization, sponsors should consider
125 during the planning phases of the study which data standards are applicable. This allows for
126 streamlining data collection (e.g., data collection instruments and processes can be designed to
127 collect study data in a standardized format or in a format that is amenable to standardization) and
128 ensuring that all relevant data elements are collected in the proper format. Some sponsors
129 choose to convert nonstandardized data to a standardized format. In general, we discourage this
130 conversion approach unless it is the only option.¹¹ These two approaches (collecting
131 standardized data vs. conversion of nonstandardized collected data) have somewhat different
132 considerations that are discussed below.

133 134 **A. Planning and Providing Standardized Study Data**

135
136 For clinical and nonclinical studies, sponsors should describe in the IND or IDE their plan to
137 submit standardized study data to the Agency. For an IND, the plan should be located in the
138 general investigational plan.¹² For an IDE, the plan should be part of the IDE study protocol
139 under the *Data Management Plan*. This study data standardization plan should contain a listing
140 of the planned studies, or a description of the types of studies, along with the data standards that
141 the sponsor intends to use. Where it is not possible or practical to use the FDA-recommended
142 standards, the study data standardization plan should describe the studies or types of studies that
143 will not be standardized, and describe why the use of FDA-recommended standards is not
144 feasible. The inclusion of a study data standardization plan gives the Agency the opportunity to
145 identify data standardization issues early in the development program. Sponsor feedback on
146 feasible use of FDA-recommended standards may provide future opportunities to refine or
147 develop data standards. The study data standardization plan should be updated in future
148 communications with the Agency as the development program evolves.

149
150 When using a data standard, there may be occasional ambiguity resulting in more than one way
151 to implement the data standard.¹³ Instances in which a standard allows for more than one
152 implementation should ideally be discussed with the appropriate review division prior to data
153 submission. Submitters should ensure their data are fully compliant with the standard (see

⁹ See <http://www.fda.gov/ForIndustry/DataStandards/StudyDataStandards/default.htm>.

¹⁰ By *supported*, we mean that the FDA has established processes and technology infrastructure to support the receipt, processing, review, and archiving of study data using these standards.

¹¹ The conversion approach after the study is complete may disclose problems with data collection that could have been corrected earlier had standardized data been collected from the outset.

¹² See 21 CFR 312.23(a)(3) for a description of the general investigational plan. Note that a study data standardization plan is not a mandatory portion of the general investigational plan, but we recommend it be included in this section of the IND.

¹³ For example, the CDISC SDTM Implementation Guide v. 3.1.2 describes the reference start date (RFSTDTC) as *usually* equivalent to date/time when subject was first exposed to study treatment. The word *usually* indicates other interpretations of the reference start date are possible.

Contains Nonbinding Recommendations

Draft – Not for Implementation

154 section III.D, Data Validation). Any data that do not easily conform to the standard should be
155 described in the study data submission.

156
157 At the time of a study data submission, submitters should describe in the cover letter how the
158 study data standardization plan was implemented.¹⁴ In addition, the submission should include a
159 tabular listing of studies that contains a description of which data standards were used, including
160 versions.¹⁵

161
162 A particular data standard does not require all data elements defined by the standard to be
163 collected in any given study. For example, the Study Data Tabulation Model (SDTM)¹⁶
164 classifies variables as required, expected, or permissible. As previously described, *what* is
165 collected and submitted is the subject of science, regulation, and discussions with the review
166 division. However, all data collected in a study should be submitted, and with the highest level
167 of standardization possible.

168
169 Standardized electronic datasets should generally be accompanied by documentation typically
170 referred to as a *Reviewer's Guide*. The Reviewer's Guide describes any special considerations or
171 directions that will facilitate a reviewer's use of the datasets and can help the reviewer understand
172 the relationships between the study report and the data. The Reviewer's Guide is an integral part
173 of a standards-compliant study data submission. No additional guidance is available regarding
174 the format and content of a Reviewer's Guide at this time.

175 176 **B. Controlled Terminologies**

177
178 The use of terminology standards, also known as controlled terminologies, is an important
179 component of standardization and is a critical component of achieving semantically interoperable
180 data exchange.¹⁷ Terminology standards specify the key concepts that are represented as
181 preferred terms, definitions, synonyms, codes, and code system. Terminology standards are
182 maintained by external organizations (i.e., external to the submitter). Sponsor-defined custom
183 terms are not considered controlled terminologies. Examples of controlled terminologies
184 include:

- 185
- 186 • MedDRA (the Medical Dictionary for Regulatory Activities) for adverse events¹⁸
 - 187 • The National Drug File — Reference Terminology for drug classifications¹⁹
 - 188 • The International Standards Organization (ISO) 3166-1 alpha-3 standard for the countries
189 of the world²⁰

¹⁴ This could include a simple statement saying that the standardization plan dated yyyy-mm-dd was executed faithfully, or that particular studies were not standardized according to the plan for stated reasons.

¹⁵ For clinical studies to CBER and CDER, this information can be included in section 5.2 Tabular Listing of All Clinical Studies within the eCTD (electronic common technical document).

¹⁶ See <http://www.cdisc.org>.

¹⁷ See the Appendix for a detailed discussion of semantic interoperability.

¹⁸ See <http://www.meddrasso.com/>.

¹⁹ See <http://ncit.nci.nih.gov/ncitbrowser/pages/vocabulary.jsf?dictionary=National%20Drug%20File%20-%20Reference%20Terminology>.

²⁰ See <http://www.fda.gov/ForIndustry/DataStandards/StructuredProductLabeling/ucm162567.htm>.

Contains Nonbinding Recommendations

Draft – Not for Implementation

- 190 • CDRH Event Problem Codes — used by CDRH to describe device problems and patient
191 problems associated with device use²¹

192

193 The analysis of study data is greatly facilitated by the use of controlled terms for clinical or
194 scientific concepts that have standard, predefined meanings and representations. The use of
195 standard terminology for adverse events is perhaps the earliest example of using data standards
196 for study data. *Myocardial infarction* and *heart attack* are both synonyms for the same clinical
197 concept, and as such should be mapped to the same term in a standard dictionary. This facilitates
198 an efficient analysis of events that are coded to the standard term. Because assigning standard
199 terms from a dictionary to actual concepts is sometimes as much an art as a science, FDA
200 expects submitters to provide in the electronic study data submission the actual verbatim terms
201 that were collected on the case report form (the so called *raw* term), as well as the coded term, so
202 that review staff can evaluate the standardization process.

203

204 The use of standard terminology is particularly useful when applied across studies, as this greatly
205 facilitates appropriate integrated analyses that are stratified by study and related cross-study
206 analyses (e.g., when greater power is needed to detect important safety signals). Cross-study
207 comparisons and multi-study pooled analyses frequently provide critical information for
208 regulatory decisions, such as confirmation of efficacy,²² exposure-response relationships,²³ and
209 population pharmacokinetics.²⁴

210

211 The FDA Study Data Standards Resources Web page²⁵ contains the currently supported
212 terminology standards at CBER, CDER, and CDRH. When planning a study (including the
213 design of case report forms, data management systems, and statistical analysis plans), the
214 sponsor should identify which FDA-supported standard terminologies to use for submission. In
215 cases in which a sponsor does not intend to use an FDA-supported standard terminology, the
216 sponsor should consult with and obtain concurrence from the appropriate review division.

217

218 If a sponsor identifies a concept for which no standard term exists, we recommend that the
219 sponsor submit the concept to the appropriate terminology maintenance organization as early as
220 possible to have a new term added to the standard dictionary. We consider this *good*
221 *terminology management practice* for any organization. The creation of custom terms for a
222 submission is discouraged (i.e., so called *extensible* code lists) as this does not support

²¹ See <http://www.fda.gov/MedicalDevices/Safety/ReportaProblem/EventProblemCodes/default.htm>

²² See the guidance for industry on *Providing Clinical Evidence of Effectiveness for Human Drugs and Biological Products*, <http://www.fda.gov/downloads/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/ucm072008.pdf>. We update guidances periodically. To make sure you have the most recent version of a guidance, check the FDA Dugs guidance Web page at <http://www.fda.gov/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/default.htm>.

²³ See the guidance for industry on *Exposure-Response Relationships — Study Design, Data Analysis, and Regulatory Applications*,

<http://www.fda.gov/downloads/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/ucm072109.pdf>.

²⁴ See the guidance for industry on *Population Pharmacokinetics*,

<http://www.fda.gov/downloads/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/ucm072137.pdf>.

²⁵ See <http://www.fda.gov/ForIndustry/DataStandards/StudyDataStandards/default.htm>.

Contains Nonbinding Recommendations

Draft – Not for Implementation

223 semantically interoperable study data exchange. Terminology maintenance organizations
224 generally have well-defined change control processes, and sponsors should allow sufficient time
225 for a proposed term to be vetted, as it is desirable to have the term incorporated into the standard
226 terminology before the data are submitted. If the use of custom terms cannot be avoided, the
227 submitter should clearly identify and define the custom terms within the submission and use
228 them consistently throughout the application. This is particularly important for primary and
229 secondary clinical endpoints and patient reported outcomes used as endpoints.

230
231 If a submitter identifies an entire information domain²⁶ for which FDA has not endorsed a
232 specific standard terminology, then the submitter may select a standard terminology to use. We
233 recommend that the submitter discuss this selection with the review division in advance. The
234 same terminology should then be used consistently throughout all relevant studies within the
235 application.

236
237 A frequent question is how to handle multiple versions of a dictionary within a single
238 application. We recognize that studies are done at different times, during which different
239 versions of a dictionary may be the most current. We expect sponsors to use the most recent
240 version of an FDA-supported dictionary available at the time of coding (recognizing that a single
241 study will use a single version).²⁷ It is understandable and therefore acceptable to have different
242 studies use different versions of the same dictionary within the same application. However,
243 important pooled analyses of coded terms across multiple studies (e.g., a pooled adverse event
244 analysis of all pivotal trials in an integrated summary of safety) should be conducted using the
245 same version of a dictionary. This is the only way to ensure a truly consistent and coherent
246 comparison of clinical and scientific concepts across multiple studies. Sponsors should clearly
247 document which terminologies and which versions of terminologies are used for every study
248 within a submission (see section III.A). Submitters should direct any specific questions
249 regarding the use of controlled terminologies to the appropriate review division.

250

C. Standardization of Previously Collected Nonstandard Data

251

252
253 Clinical and nonclinical study data that were previously collected in a nonstandard format are not
254 always easily amenable to complete standardization. Typically, a conversion to a standard
255 format will map every data element as originally collected to a corresponding data element
256 described in a standard. Some study data conversions will be straightforward and will result in
257 complete data in a standardized format. In some instances, it may not be possible to represent a
258 collected data element as a standard data element. In these cases, the submission should
259 document why certain data elements could not be fully standardized or were otherwise not
260 included in the standardized data submission. In cases where the data were collected on a CRF
261 but not included in the converted datasets, the omitted data should be apparent on the annotated
262 CRF. The tabular listing of studies in a submission should indicate which studies contained
263 previously collected nonstandard data that were subsequently converted to standard format.

²⁶ By *information domain*, we mean a logical grouping of clinical or scientific concepts that are amenable to standardization (e.g., adverse event data, laboratory data, histopathology data, imaging data).

²⁷ That is to say, if a new version of a dictionary is released in the middle of the coding process for a given study, then the sponsor should continue to use the version that was most current when the coding process began.

Contains Nonbinding Recommendations

Draft – Not for Implementation

264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303

D. Data Validation

For purposes of this guidance, data validation is a process that attempts to ensure that data are both compliant and useful. *Compliant* means the data conform to the applicable data standards. *Useful* means the ability of the data to support the intended use (e.g., regulatory review and analysis). Data validation is one method used to assess data quality.²⁸ Data validation relies on a set of validation rules that are used to check that the data conform to a minimum set of quality standards. Data validation attempts to identify early in the review common or serious data issues that may adversely affect use of the data. We recognize that data validation is an imperfect process, as it is impossible or impractical to define *a priori* all the relevant validation rules for any given submission. Sometimes serious deficiencies in the data are only evident through manual inspection of the data and may only become evident once the review is well under way. Often these are because of problems in data content (i.e., *what* was or was not submitted, or issues with the original source data), and not necessarily in *how* the data were standardized. We expect the data validation rules will undergo improvements as the Agency incrementally refines its definition of valid, high quality data for regulatory decisions.

Cross-study data validation for consistency is especially important for analyses supporting the confirmation of efficacy, exposure-response relationships, and population pharmacokinetics. For example, the terminology of the observational data elements (see section III.B) and the units of their values should be consistent across relevant studies.

FDA generally recognizes two types of validation rules:

Technical validation rules help ensure that the data conform to the data standards. For example, a technical validation rule for Clinical Data Interchange Standards Consortium (CDISC) Study Data Tabulation Model (SDTM) data would check that the value in the DOMAIN column of all datasets matches the name of the domain.

Business validation rules help ensure that the data will support business processes that rely on the data (i.e., support meaningful analysis of the data). For example, a business validation rule for a human study may require that each value for AGE fall in within a prespecified human physiologic range.

Once a data standard is defined, the technical validation rules are rather static. They are not expected to change substantially unless the standard itself changes. Business validation rules may evolve over time as new analysis requirements are identified and incorporated into data validation processes.

²⁸ Data quality is beyond the scope of this guidance. For purposes of this guidance, standardized data do not ensure quality data; however, standardized data do make it easier to assess some aspects of data quality by facilitating the automation of various data quality checks (e.g., completeness, reasonableness). Data standards can also improve the collection of quality data by facilitating the use of standard data quality checks in electronic data collection instruments.

Contains Nonbinding Recommendations

Draft – Not for Implementation

304 Submitters can find the validation rules used at FDA on the Study Data Standards Resources
305 Web page.²⁹ Submitters should validate their data before submission using the published
306 validation rules and either correct any validation errors or explain in the submission why certain
307 validation errors could not be corrected. The submission should state the version of the
308 validation rules used.

309
310 Upon submission, FDA will conduct its own data validation, using the same version of the
311 validation rules used by the sponsor. FDA will use the data validation results to inform review
312 staff of potential problems in using the data, and to assess the usefulness of the rules. If
313 necessary, FDA will report serious data validation errors to the submitter for correction. Please
314 refer to the online validation rules for a description of data validation errors. The recommended
315 presubmission validation step is intended to minimize the presence of serious validation errors at
316 the time of submission.

E. Exceptions to Standardized Study Data Submissions

317
318
319
320 Data standards are evolving at different rates for varying study types, and they are not being
321 adopted at the same rate throughout FDA. Accordingly, there may be some exceptions to FDA's
322 general expectation that all study data be standardized for submission. The data standards
323 catalog on the Study Data Standards Resources Web page describes which types of studies
324 should be submitted in a standardized format, and which offices are currently accepting
325 standardized study data submissions.

F. Meetings With FDA

326
327
328
329 Sponsors can use established FDA-sponsor meetings (e.g., pre-IND, pre-IDE, and end-of-phase
330 2) to raise data standardization issues (if any). Discussions about nonclinical study data
331 standardization plans can be initiated at the pre-IND stage and should continue throughout
332 development. Initial discussions about which data standards to use for clinical study data should
333 take place as early as possible during drug development, especially for safety data, which is often
334 less well planned for in advance, but should occur no later than the end of phase 2. In general,
335 the premarketing application meeting is considered too late to initiate data standardization
336 discussions.

337
338 Sponsors and applicants may submit technical questions related to data standards at any time to
339 the technical support team identified by each Center (see the Study Data Standards Resources
340 Web page for specific contact information). Submitters may also request a separate Type C
341 meeting to discuss substantive data standardization issues. An example of such an issue might
342 be a sponsor's desire to use a non-supported terminology. The request should include adequate
343 information to identify the appropriate FDA staff necessary to discuss the proposed agenda
344 items.

345
346

²⁹ See <http://www.fda.gov/ForIndustry/DataStandards/StudyDataStandards/default.htm>.

Contains Nonbinding Recommendations

Draft – Not for Implementation

347 **IV. ONLINE TECHNICAL RESOURCES**

348
349 The FDA maintains online study data standards technical resources at
350 <http://www.fda.gov/ForIndustry/DataStandards/StudyDataStandards/default.htm>, which includes
351 the following:

- 352
- 353 • A data standards catalog that lists the various data standards (including versions of those
- 354 standards) currently accepted at each Center
- 355 • Study data technical specifications document
- 356 • Study data validation rules
- 357 • Other Center-specific resources
- 358

359 Submitters should use the data standards listed on this Web page for their study data
360 submissions. Data standards for study data can be expected to evolve over time; therefore,
361 submitters are encouraged to visit the Web page regularly to obtain the latest information.

362
363 Each Center maintains a Data Standards Plan that outlines the management and adoption of data
364 standards. As new standards are supported, the Agency will support older standards and versions
365 of standards for at least 2 years after the addition of a new standard or version. FDA will update
366 the Web page as needed to reflect changes in data standards and associated technical
367 specifications.

368 369 370 **V. ADDITIONAL SUPPORT**

371
372 Sponsors/applicants with questions on how to implement the FDA-supported study data
373 standards should contact and work with FDA technical staff. Contact information is provided on
374 the Study Data Standards Resources Web page.³⁰ Sponsors may also arrange to submit sample
375 data for a presubmission technical review. The technical staff also welcomes any additional
376 feedback or comments regarding the information posted on the Web page.

³⁰ See <http://www.fda.gov/ForIndustry/DataStandards/StudyDataStandards/default.htm>.

Contains Nonbinding Recommendations

Draft – Not for Implementation

377 **APPENDIX – DATA STANDARDS AND INTEROPERABLE DATA EXCHANGE**

378

379 This appendix provides information on why FDA is now emphasizing the submission of
380 standardized data for both nonclinical and clinical studies. It also provides some of the guiding
381 principles for the Agency’s long-term study data standards management strategies. An important
382 goal of standardizing study data submissions is to achieve an acceptable degree of *semantic*
383 *interoperability* (discussed below). This appendix describes different types of interoperability
384 and how data standards can support interoperable data exchange now and in the future.

385

386 At the most fundamental level, study data can be considered a collection of data elements and
387 their relationships. A data element is the smallest (or *atomic*) piece of information that is useful
388 for analysis (e.g., a systolic blood pressure measurement, a lab test result, a response to a
389 question on a questionnaire).

390

391 A data value is by itself meaningless without additional information about the data (so called
392 *metadata*). Metadata is often described as *data about data*. Metadata is structured information
393 that describes, explains, or otherwise makes it easier to retrieve, use, or manage data.³¹ For
394 example, the number *44* itself is meaningless without an association with *Hematocrit*.
395 Hematocrit in this example is metadata that further describes the data.

396

397 Just as it is important to standardize the representation of data (e.g., M and F for male and
398 female, respectively), it is equally important to standardize the metadata. The expressions
399 Hematocrit = 44; Hct = 44, or Hct Lab Test = 44 all convey the same information to a human,
400 but an information system or analysis program will fail to recognize that they are equivalent
401 because the metadata is not standardized. It is also important to standardize the definition of the
402 metadata, so that the meaning of a Hematocrit is constant across studies and submissions.

403

404 In addition to standardizing the data and metadata, it is important to capture and represent
405 relationships (also called associations) between data elements in a standard way. Relationships
406 between data elements are critical to understand or interpret the data. Consider the following
407 information collected on the same day for one subject in a study:

³¹ Metadata is said to “give meaning to data” or to put data “in context.” Although the term is now frequently used to refer to XML (extensible markup language) tags, there is nothing new about the concept of metadata. Data about a library book such as author, type of book, and the Library of Congress number, are metadata and were once maintained on index cards. SAS labels and formats are a rudimentary form of metadata, although they have not historically been referred to as metadata. Additional discussion about metadata is outside the scope of this guidance other than to stress the importance of standardized metadata for study data submissions.

Contains Nonbinding Recommendations

Draft – Not for Implementation

408
409 Systolic Blood Pressure = 90 mmHg
410 Position = standing
411 Systolic Blood Pressure = 110 mmHg
412 Time = 10:23 a.m.
413 Time = 10:20 a.m.
414 Position = lying

415
416 When presented as a series of unrelated data elements, they cannot reliably be interpreted. Once
417 the relationships are captured, as shown simply below using arrows:

418 Time = 10:20 a.m. \leftrightarrow Position = lying \leftrightarrow Systolic Blood Pressure = 110 mmHg
419 Time = 10:23 a.m. \leftrightarrow Position = standing \leftrightarrow Systolic Blood Pressure = 90 mmHg
420

421
422 the interpretation of a drop in systolic blood pressure of 20 mmHg while standing, and therefore
423 the presence of clinical orthostatic hypotension, is possible. Standardizing study data therefore
424 involves standardizing the data, metadata, and the representation of relationships.

425
426 With these fundamental concepts of data standardization in mind, data standards can be
427 considered in the context of interoperable data exchange.

428 429 **A. Interoperability**

430
431 Much has been written about interoperability, with many available definitions and interpretations
432 within the health informatics community. In August 2006, the President signed an Executive
433 Order mandating the Federal Government use of interoperable data standards for health
434 information exchange.³² Although this order was directed at Federal agencies that administer
435 health care programs (and therefore not FDA), it is relevant to this guidance because it defined
436 interoperability:

437
438 *“Interoperability” means the ability to communicate and exchange data accurately, effectively,*
439 *securely, and consistently with different information technology systems, software applications,*
440 *and networks in various settings, and exchange data such that clinical or operational purpose*
441 *and meaning of the data are preserved and unaltered.*

442
443 Achieving interoperable study data exchange between submitters and FDA is not an all-or-
444 nothing proposition. Interoperability represents a continuum, with higher degrees of data
445 standardization resulting in greater interoperability, which in turn makes the data more useful
446 and increasingly capable of supporting efficient processes and analyses by the data recipient. It
447 is therefore useful to understand the degree of interoperability that is desirable for standardized
448 study data submissions.

449

³² See <http://www.cga.ct.gov/2006/rpt/2006-R-0603.htm>.

Contains Nonbinding Recommendations

Draft – Not for Implementation

450 In 2007, the Electronic Health Record Interoperability Work Group within Health Level Seven
451 issued a white paper that characterized the different types of interoperability based on an analysis
452 of how the term was being defined and used in actual practice.³³ Three hierarchical types of
453 interoperability were identified: technical, semantic, and process interoperability. A review of
454 these three types provides insight into the desired level of interoperability for standardized study
455 data submissions.

456
457 **Technical interoperability** describes the lowest level of interoperability whereby two different
458 systems or organizations exchange data so that the data are useful. There is nothing that defines
459 how useful. The focus of technical interoperability is on the conveyance of data, not on its
460 meaning. Technical interoperability supports the exchange of information that can be used by a
461 person but not necessarily processed further. When applied to study data, a simple exchange of
462 nonstandardized data using an agreed-upon file format for data exchange (e.g., SAS transport
463 file) is an example of technical interoperability.

464
465 **Semantic interoperability** describes the ability of information shared by systems to be
466 understood, so that nonnumeric data can be processed by the receiving system. Semantic
467 interoperability is a multi-level concept with the degree of semantic interoperability dependent
468 on the level of agreement on data content terminology and other factors. With greater degrees of
469 semantic interoperability, less human manual processing is required, thereby decreasing errors
470 and inefficiencies in data analysis. The use of controlled terminologies and consistently defined
471 metadata support semantic interoperability.

472
473 **Process interoperability** is an emerging concept that has been identified as a requirement for
474 successful system implementation into actual work settings. Simply put, it involves the ability of
475 a system to provide the right data to the right entity at the right point in a business process.

476
477 In a regulatory setting, an example of process interoperability is the ability to quickly and
478 automatically identify and provide all the necessary information to produce an expedited adverse
479 event report in a clinical trial upon the occurrence of a serious and unexpected adverse event.
480 The timely submission of this information is required by regulation to support FDA's mandate to
481 safeguard patient safety during a clinical trial. Process interoperability becomes important when
482 particular data are necessary to support time-dependent processes.

483
484 Because the vast majority of study data are submitted after the study is complete, achieving
485 process interoperability for study data submissions in a regulatory setting is relatively
486 unimportant, at least for the foreseeable future. It is reasonable to conclude that it is most
487 desirable to achieve *semantic interoperability* in standardized study data submissions.

³³ See Coming to Terms: Scoping Interoperability for Health Care <http://www.hln.com/assets/pdf/Coming-to-Terms-February-2007.pdf>.

Contains Nonbinding Recommendations

Draft – Not for Implementation

488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532

B. Types of Data Standards

There are various types of data standards, and all have a role in achieving semantically interoperable study data exchange. The following describe working definitions that are not necessarily mutually exclusive.

A **file format standard** specifies a particular way that information is encoded in a computer file. Specifications for a format permit the file to be written according to a standard, opened for use or alteration, and written back to a storage medium for later access. Some file formats in widespread use are proprietary, others are open source. Examples of file format standards supported at FDA include Adobe Acrobat Portable Document (.pdf), SAS Transport File format (.xpt), text files (.txt), and Extensible Markup Language (.xml). The use of a file format standard for study data exchange supports technical interoperability, but by itself is often insufficient for semantic interoperability.

An **exchange standard** describes a standard way of exchanging data between computer systems. Exchange standards may describe the data elements and relationships necessary to achieve the unambiguous exchange of information between disparate information systems. Examples include the Health Level 7 (HL7) Structured Product Labeling (SPL), Individual Case Safety Report (ICSR), and ECG Waveform standards. The use of an exchange standard is necessary but often not sufficient to achieve semantically interoperable data exchange. Often the exchange standard should be combined with standard terminologies to achieve an acceptable degree of semantic interoperability.

A **terminology standard** is the group of specifications for commonly agreed-upon vocabulary to be used in another data standard or family of data standards. Terminology standards specify the key concepts that are represented as preferred terms, definitions, synonyms, codes, and the code system that must be used to ensure a common understanding of a given data standard. Examples include the ISO 3166-1 alpha-3 standard for representing the countries of the world, the FDA/USP Substance Registration System’s Unique Ingredient Identifier (UNII) for representing chemical and biological substances, the Department of Veterans Affairs’ National Drug File – Reference Terminology for classifying drugs. Terminology standards, when used within exchange and file format standards, combine to provide a useful degree of semantically interoperable data exchange.

An **analysis standard** describes a standard presentation of the data intended to support analysis. It includes extraction, transformation, and derivations of the original data. An example is the CDISC Analysis Data Model (ADaM) standard.

As a practical matter, most of the existing data standards do not have one-to-one correspondence with these types of data standards. For example, the CDISC Study Data Tabulation Model (SDTM) can be used as an exchange standard (when combined with the SAS transport file format standard) and as an analysis standard to support simple analyses, such as demographics or adverse events.

Contains Nonbinding Recommendations

Draft – Not for Implementation

533
534 Those familiar with data standards often describe other types of standards: data organization
535 standards, content standards (such as disease/domain specific standards (e.g., for Alzheimer’s
536 disease or Diabetes Mellitus)), and data presentation standards. We find these are all variations
537 on the previously described data standard types.
538
539 In summary, the goal of standardizing study data is to make the data more useful and to support
540 semantically interoperable data exchange between submitters and the FDA. This, in turn, will
541 support more efficient analytic processes at the Agency. The combined use of various types of
542 data standards is necessary to achieve an acceptable degree of semantic interoperability.