

Evolution and Implementation of the CDISC Study Data Tabulation Model (SDTM)

Tom Guinter

VP Clinical Data Strategies, Octagon Research Solutions, and CDISC

Fred Wood

Global Standards Manager, Procter and Gamble Pharmaceuticals, and CDISC

ABSTRACT

The CDISC (Clinical Data Interchange Standards Consortium) SDTM is a standard for submitting data tabulations to the FDA in support of marketing applications. In July of 2004, this standard became Study Data Specification (1) referenced in the eCTD Guidance (2). This paper/presentation will provide an overview of the SDTM and the associated Implementation Guides, commonly referred to as the SDTMIG and SEND. Included will be the evolution of the standard, the current status of the SDTM, strategic reasons for considering implementing the SDTM, and how the SDTM is used in the FDA review environment, including an overview of the tools used by the FDA to review SDTM data.

INTRODUCTION

CDISC BACKGROUND

From its inception in 1997, CDISC has recognized the need for the establishment of standard data models to improve the process of electronic acquisition and exchange of clinical trials information for the benefit of all medical and pharmaceutical stakeholders. This is reflected in the mission statement: *"CDISC is an open, multidisciplinary, non-profit organization committed to the development of worldwide industry standards to support the electronic acquisition, exchange, submission and archiving of clinical trials data and metadata for medical and biopharmaceutical product development. The mission of CDISC is to develop and support global, platform-independent data standards that enable information system interoperability to improve medical research and related areas of healthcare."*

CDISC has four major standards: the Lab Model (LAB) for the transfer of lab data from vendors to sponsors; the Operational Data Model (ODM) for the transfer of clinical trial data and metadata, including administrative data and audit trail; the Analysis Data Model for analysis datasets; and the SDTM for data tabulations. The LAB and ODM models are designed for the transfer of data from vendors and CROs to sponsors, while the ADaM and SDTM models apply to the submission of data to the FDA.

CRTs

Case Report Tabulations originated in FDA regulation of the late 1980's. Originally intended to allow sponsors to submit data tabulations to the FDA instead of copies of every subject CRF. Recently, as documented in the FDA Study Data Specification (1), which is referenced in the eCTD Guidance (2), the definition has expanded to include data listings, patient profiles, data tabulations, and analysis datasets.

EVOLUTION OF THE SDTM

The SDTM originated as the Submission Data Model (SDM), developed by the CDISC Submission Data Standards (SDS) Team. The SDS Team is comprised of approximately thirty volunteers from industry. Team meetings have also been regularly attended by three key FDA observers representing the Office of Business Process Support, Statistics, and Medical Review. Team members come from a cross-section of industry, including most major pharmaceutical companies, numerous small to mid-level pharmaceutical companies, CROs, and service providers. The SDS Team collaborates regularly through biweekly teleconferences, quarterly face-to-face meetings, and literally thousands of regular email communications, all in an effort to support the CDISC mission, to "...improve medical research and related areas of healthcare."

The concept that developed into the Study Data Tabulation Model (SDTM) v1.0 / Submission Data Standards (SDS) v3.1 was initially presented to the SDS Team by FDA liaisons in October, 2002. Prior to that time, the SDS Team had developed the Submission Data Domain Models Version 1.0 (v1, 2001), and Version 2.0 (v2, early 2002), and was just about to publish v2.1. All of these versions focused on safety data/domains. While the v1 and v2 safety data concepts had been well received in industry, it was recognized there were a couple of major shortcomings. The revolutionary concept that the FDA proposed, termed Version 3.0 by the Team, addressed those major shortcomings by 1) providing a standard for all clinical trial data, not

just safety data, and 2) providing a standard based on data modeling principles, rather than data management / operational database principles. This provided a consistent approach for modeling all clinical trial data across three primary data types (Events, Interventions, and Findings), allowing FDA review tools to be developed.

The preliminary draft version of the SDTM concept was published in June, 2003 as the Submission Data Domain Models Version 3.0, or better known as SDS v3.0. The first version intended for implementation was published as two documents in June, 2004: the SDTM v1.0 (the model), and SDS v3.1 (the implementation guide). During the period from June 2003 to June 2004, there were a number of enhancements leading to the final approved version, many the result of a joint FDA/industry pilot to test SDS v3.0, and two public review/comment periods. Shortly after the publication of the production version in June, 2004, the FDA recognized the SDTM as an approved method for submitting the data-tabulation component of Case Report Tabulations (CRTs) in the then-draft eCTD guidance. Health and Human Services (HHS) announced in the December, 2004 and May, 2005 regulatory agendas that the agency was moving towards requiring the submission of clinical trial data in the SDTM format. While we, of course, understand regulatory changes take time, this clearly shows the FDA's commitment to the SDTM.

Dr. Janet Woodcock, acting Director of CDER, met with the CDISC Board of Directors (BOD) in January, 2005 to review the FDA's Critical Path Initiative. She clearly stated that industry adoption of the SDTM for submission of data for all clinical trials in marketing applications was a significant component of the FDA's Critical Path Initiative, and asked the CDISC BOD for recommendations on how CDISC and the FDA could enhance their collaboration to promote quicker industry adoption of the SDTM.

On February 1, 2005, the FDA conducted a public meeting to review the status of industry adoption of the SDTM. Multiple FDA liaisons reviewed the importance of the SDTM and its value to both industry and FDA. They also confirmed that sponsors submitting data in the SDTM would not be required to submit data listings and patient profiles. FDA review tools will automatically produce data listings and patient profiles from properly formatted SDTM datasets. Additionally, FDA agreed to provide advance training to reviewers on the SDTM, and the review tools, when sponsors notify FDA in advance that they are submitting in the SDTM format. The FDA has recognized that ensuring this training occurs is a necessary component of SDTM implementation in industry.

Beyond the formal communications, we have seen considerable effort invested by multiple FDA industry liaisons, in numerous presentations at industry, promoting the value of the SDTM to the FDA.

LONG-TERM BENEFITS

The FDA has been developing the Janus data warehouse as the repository to store all the submitted data. It will provide a stable foundation for FDA's growing list of standard review tools, developed through Cooperative Research and Development Agreements (CRADAs) between the agency and various vendors. The tools have been developed to utilize the SDTM, rather than vice versa. Most of these tools create tabular and/or graphical views of the SDTM data via canned reports, although custom reports can be run as well. The tools routinely bring demographics and treatment data into all views of subject data, utilizing the fixed relational database structure of Janus. Many tools also provide built-in hyperlinks that allow drilling down from group summaries to individual-subject data, and to navigate from graphs to tables. Some allow for various "what-if" scenarios, such as allowing the elimination of outlying values from the calculation of means.

From a strategic standpoint, it is to a sponsor's advantage to move toward implementation of the SDTM for submission datasets. Once FDA reviewers become accustomed to being able to easily navigate through submitted data using the dedicated tools, being unable to do so for a subsequent submission may lead to inefficiencies in the review process.

SDTM BASICS

The current CDISC Submission Data Standard consists of two documents: the Study Data Tabulation Model Version 1.1, published in April 2005, and the Study Data Tabulation Model Implementation Guide: Human Clinical Trials Version 3.1.1, published in August 2005. The first describes the model, while the second provides guidance on model implementation, including domain models and examples with real data for commonly submitted datasets, a set of assumptions to aid in interpretation of the intended implementation, more detailed descriptions of the Trial Design Model (TDM) tables, and a more detailed discussion on representing relationships within and across submission datasets. The SDTM is built around several key concepts. These are described in the following paragraphs.

Domains

Domains are groups of related observations. These observations are grouped by topic in datasets. Datasets and domains are usually the same, but some domains contain two classes of observations and have to be split into two datasets.

Observations

Observations can be described by a series of named variables. Each variable, which normally corresponds to a column in a dataset, can be classified according to its role. Most variables in a domain begin with a prescribed prefix.

Example: In Study ABC001, Subject 1234-0001 had a heart rate of 100 bpm on Study Day 6. This would be represented in a dataset as follows:

STUDYID	USUBJID	VSDY	VSTESTCD	VSORRES	VSORRESU
ABC001	1234-0001	6	HR	100	bpm

Observation Classes

An observation can be classified as one of three major types: Interventions, Events, or Findings.

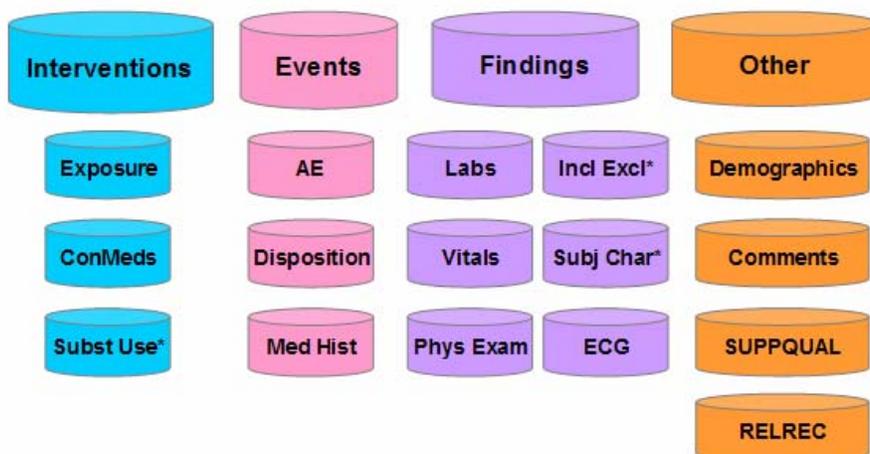
Interventions: investigational treatments, therapeutic treatments, and surgical procedures administered to the subject or animal. One record per constant dosing/treatment interval.

Events: occurrences or incidents independent of planned study evaluations occurring during the trial (e.g., adverse events) or prior to the trial (e.g., medical history). One record per event.

Findings: observations resulting from planned evaluations (e.g. lab tests, ECGs, microscopic findings). One record per finding result or measurement.

Each observation class has a defined set of standard variables. Only variables defined for an observation class can be used in the respective SDTM dataset. Other variables for which a sponsor collected data must be represented in a SUPPQUAL dataset, to be discussed later. The standard variables, coupled with identifiers and timing variables which are used in all observation classes are the building blocks for constructing SDTM domains. Figure 1 shows the observation class for some commonly collected domains.

Figure 1. Fitting V3.1 Domains into Observation Classes



Variable Roles

Every variable has been assigned a Role which describes the type of information conveyed by each variable within an observation. These Roles are defined in the SDTM:

- **Topic Variable** - Identifies the focus of the observation. There is only one per dataset.
- **Identifier Variables** - Identify the study, the subject, the domain, and sequence number of the observation.
- **Timing Variables** - Describe the start and end of the observation, and/or when it was collected.
- **Qualifier Variables** - Describe the attributes and results of the observation. These can be further subdivided into Grouping, Result, Synonym, Record, and Variable Qualifiers.

Variable Metadata

Each dataset or table is accompanied by metadata definitions that provide information about the variables used in the dataset. Included are the SAS label, the {data} type, controlled terms or format, the origin (e.g., CRF, derived), and the role (as described above).

Additional Datasets and Tables which Need To Be Submitted

In addition to data submitted in accordance with the three observation classes, there are a number of special-purpose datasets which are also part of the SDTM. Included are the Demographics dataset, the Comments dataset, the SUPPQUAL dataset(s), the RELREC (Related Records) dataset, and seven TDM datasets.

The Demographics dataset includes a set of standard variables that describe each subject in a clinical study. The Comments domain is a fixed domain that provides a solution for submitting free-text comments related to data in one or more domains or collected on a separate CRF page dedicated to comments. Comments is similar to the Supplemental Qualifiers (SUPPQUAL) dataset but it allows for one comment to span multiple variables (COVAL-COVALn) in order to accommodate comments longer than 200 characters. Comments may be related to a Subject or to specific parent records in an SDTM domain.

Relationship Tables

In order to understand the two primary relationship tables, SUPPQUAL and RELREC, it is necessary to understand that every record in the observation-class datasets has a unique set of keys consisting of STUDYID, USUBJID, DOMAIN, and --SEQ (the two hyphens indicating the two-letter domain code). The --SEQ variable is sponsor defined, and unique within SUBJECT, DOMAIN, and STUDY. Knowing the values for these four columns allows precise identification of a record within a submission. Both RELREC and SUPPQUAL use this same concept to identify a parent record to which either a non-standard variable (SUPPQUAL) or another independent (parent) record (RELREC) is related. These tables use STUDYID, USUBJID, RDOMAIN (related domain), and two other variables which point to the related record. These two variables are IDVAR and IDVARVAL, which describe the parent record's unique key variable and value.

The Supplemental Qualifiers (SUPPQUAL) special-purpose dataset is used to submit values for variables not presently included in the general-observation-class models. In addition to the keys described in the previous paragraph, each SUPPQUAL record also includes the name of the Qualifier variable being added (QNAM), the label for the variable (QLABEL), the actual value for each instance or record (QVAL), the origin (QORIG) of the value (whether it was collected via CRF, assigned or derived), and the Evaluator (QEVAL, to specify the role of the individual who assigned the value, such as an Adjudication Committee or the sponsor).

A common reason for using a SUPPQUAL dataset is to capture attributions. An attribution is typically an interpretation or subjective classification of one or more observations by a specific evaluator, such as a population flag that classifies a subject or their data according to their evaluability for efficacy analysis. Since it is possible that different attributions may be necessary in some cases, SUPPQUAL provides a mechanism for incorporating as many attributions as are necessary. For example, if two individuals provide a determination on whether an adverse event is treatment emergent (e.g., the investigator and an independent adjudicator) then separate QNAM values should be used for each set of information, perhaps AETRTEM and AETRTEM2. This is necessary to ensure that reviewers can join/merge/transpose the information back with the records in the original domain without risk of losing information.

Many sponsors will also need to use a SUPPQUAL dataset to submit additional non-standard variables that cannot be represented in the general classes. The optional grouping identifier variable --GRPID can be a more efficient method of representing relationships in SUPPQUAL to identify individual qualifier values (SUPPQUAL records) related to multiple domain records that could be grouped, such as relating an attribution to a group of ECG measurements. If the SUPPQUAL dataset becomes too large, a sponsor has the option of submitting a separate supplemental qualifier dataset for each submitted domain that has supplemental qualifiers. The naming convention for these is "supp" followed by the two-letter domain code (e.g., suppae.xpt for adverse events).

The Related Records (RELREC) dataset is used to describe relationships between records in two (or more) datasets, such as an Event record and an Intervention record, or a Finding record and an Event record. Such a relationship is defined by creating RELREC records for each of the related observation records, and by assigning a unique character identifier value for the relationship. Each RELREC record contains same keys as SUPPQUAL to identify a record (using --SEQ in IDVAR) or group of records (using --GRPID in IDVAR). RELREC uses an additional, unique variable, RELID, the relationship identifier, which is the same for all related records. The value of RELID can be any constant value chosen by the sponsor.

The Trial Design Model (TDM)

The TDM allows description of key aspects of the planned conduct of a clinical trial in a standardized way. These standardized descriptions will allow reviewers to:

- clearly and quickly grasp the design of a clinical trial
- compare the designs of different trials
- search a data warehouse for clinical trials with certain features
- compare planned and actual treatments and visits for subjects in a clinical trial

Modeling a clinical trial in this standardized way requires the explicit statement of certain decision rules that may not be addressed, or may remain vague or ambiguous, in the usual prose protocol document. Prospective modeling of the design of a clinical trial should lead to a clearer, better protocol. Retrospective modeling of the design of a clinical trial should ensure a clear description of how the trial was interpreted by the sponsor.

The TDM is built upon the concepts of Elements, Arms, Epochs, and Visits. They can be described as follows:

- An Element is the basic building block for time within a clinical trial, and has the following characteristics: a description of what treatment(s) are planned for the subject during the Element, a definition of the start of the Element, a rule for ending the Element.
- An Arm is a planned sequence of Elements, typically equivalent to a treatment group. Branches may take place between one Element and the next, and some designs allow for some flexibility of Elements within an Arm.
- The term EPOCH is used to refer to the portion of a blinded trial that usually corresponds to the time period for an individual Element.
- A Visit is defined as a clinical encounter that encompasses planned and unplanned trial interventions, procedures, and assessments that may be performed on a subject. A Visit has a start and an end, each described with a rule. A Visit need not be nested within an Element. In other words, it may start in one Element and end in another. In most blinded trials, the timing of Visits is the same for all subjects, regardless of the Arm to which they have been assigned. In these cases, the Arm is not needed to describe the timing of Visits, and is left blank in the Trial Visits domain. If the timing of Visits depends on Arm, then the complete set of Visits for each Arm should be represented in the domain.

The TDM also includes the Trial Inclusion/Exclusion dataset to describe the inclusion/exclusion criteria used to screen subjects. In contrast, the IE domain (subject-specific inclusion/exclusion criteria not met) contains the actual exceptions to those criteria for enrolled subjects.

The above information about the trial is provided in two types of tables: the trial-level tables which describe what was planned, and the subject-level tables which describe what actually happened. The trial-level tables consist of Trial Elements, Trial Arms, Trial Visits, and Trial Inclusion/Exclusion. The subject-level tables consist of Subject Elements and Subject Visits.

STANDARD FOR THE EXCHANGE OF NONCLINICAL DATA (SEND)

SEND is an implementation of the SDTM for nonclinical studies. The work on this standard began in July 2002, with an FDA pilot project announced in January, 2003. The pilot tested Version 1.0 for data from acute, subchronic, and carcinogenicity studies. Input from the pilot and efforts to more closely align this implementation with that for human clinical trials resulted in the Version 2.x standards, with current version being 2.3, posted in November 2005. As expected, some SEND domains are identical to those described in the SDTMIG for human clinical trials (e.g., ECG, Vital Signs); however, there are a number of domains specific to nonclinical research such as microscopic findings, and food and water consumption. Because the majority of animal studies are parallel studies with single treatments or single combination treatments per group with very few deviations from what was planned, SEND does not use the TDM, instead using a the implementation-specific Group Characteristics domain. A separate implementation guide is planned for reproductive toxicology studies because of the staggered timing of the phases of gestation and weaning within treatment groups, as well as the complex relationships that need to be maintained between mating partners and between parents and offspring, possibly through multiple generations.

WHAT DOES THE FUTURE LOOK LIKE?

Efforts are well underway to harmonize the SDTM with other CDISC and Health Level 7 (HL7) initiatives. The TDM in the SDTM is being harmonized with the HL7 Protocol Representation Group. The CDISC define.xml standard is now a Study Data Specification (1) referenced in the eCTD Guidance (2). It is a replacement for the traditional define.pdf, allowing much greater flexibility in the metadata describing the submitted data. The CDISC Terminology Team has been redesigned to work across all CDISC teams to develop standardized, controlled terminology across all the models. The CDISC Analysis Data Modeling (ADaM) team has developed a number of analysis-level standards, using the SDTM as the foundation. A sub-team has formed with representatives from the ADaM and SDS Teams to jointly propose a best-practice approach to producing and submitting analysis datasets and/or analysis logic, and to propose standard analysis models harmonized with the SDTM. A sub-team has also formed with representatives from the SDS and ODM teams to map the SDTM into ODM. And the list goes on. What does it all mean? CDISC and HL7 teams are working closely together towards a common endpoint, standardization and harmonization of all regulated clinical and preclinical research data.

CONCLUSION

Since becoming the standard for submission of clinical and preclinical trial data to the FDA in marketing applications, the SDTM has begun to be used by sponsors for their upstream processing to support Clinical Study Reports and Integrated Summaries. This paper covers the basics and background of the SDTM. Material presented in our session will expand on information presented here and provide additional justifications and implementation examples.

REFERENCES

1. Study Data Specifications. Version 1.1. 2005. Available via <http://www.fda.gov/cder/regulatory/ersr/ectd.htm>.

2. Guidance for Industry: Providing Regulatory Submissions in Electronic Format -- Human Pharmaceutical Product Applications and Related Submissions using the eCTD Specifications. U.S. Department of Health and Human Services, Food and Drug Administration, Center for Drug Evaluation and Research, Center for Biologics Evaluation and Research. October 2005. Available via <http://www.fda.gov/cder/regulatory/ersr/ectd.htm>.

ACKNOWLEDGEMENTS

The authors would like to thank the following:

- CDISC, their staff, and members for their commitment to creating industry data standards *that enable information system interoperability to improve medical research and related areas of healthcare*.
- Members of the CDISC SDS Team for their countless hours invested in bringing the SDTM and Implementation Guide to realization.
- FDA Liaisons to the CDISC SDS Team for their continued support.

CONTACT INFORMATION

Thomas Guinter
Vice President, Clinical Data Strategies
Octagon Research Solutions, Inc.
585 East Swedesford Road, Suite 200
Wayne, PA 19087
610 535 6500 x635
tguinter@octagonresearch.com
www.octagonresearch.com

Fred Wood
Global Data Standards Manager
Clinical Data and Information Management
Procter & Gamble Pharmaceuticals
8700 Mason-Montgomery Road
Mason, OH 45040
wood.fe@pg.com