# A Strategy for Creating Integrated Analysis Data Sets Based on the ADaM Model Using SDTM Compliant and Non-Compliant Data

Melvin S. Munsaka

Takeda Global Research & Development, Inc, One Takeda Parkway, Deerfield, IL 60015

## Abstract

A New Drug Application (NDA) typically requires an integrated database from several protocols or studies that are included as part of the submission. The Integrated Summary of Safety is in general an essential component in the review of a submission. These types of analysis differ from study level analysis primarily due to large amount of data that each study had generated prior to preparation of the Integrated Summary of Safety (ISS) or Integrated Analysis of Safety (IAS). The creation of analysis datasets to support ISS work is often very challenging because study data need to be converted and harmonized to the same format before initiation of programming and analysis. For the ISS, the specific analyses are usually performed based on predefined groupings (pooling) of studies with common elements, such as common patient population or common study designs, for example, short and long term, controlled versus uncontrolled studies, titration versus fixed dose, randomized versus non-randomized, and so on. Differences between studies designs, treatment regimens, duration of exposure, and patient populations can create barriers to data integration and can also lead to challenges in reliable interpretation and conclusions of the accumulated information. This paper will discuss a strategy to create an integrated database for safety analysis that is ADaM driven, using both SDTM and company standards as source data to the ISS database.

**Key Words:** Analysis Data Model (ADaM), Integrated Summary of Safety (ISS), Integrated Analysis of Safety (IAS), Clinical Data Interchange Standards Consortium (CDISC), Study Data Tabulation Model (SDTM), Statistical Analysis Plan (SAP), Common Technical Document (CTD)

## 1. Introduction

In CFR 21 314.50 (d) (5) (vi) (a), it is stated that: *The applicant shall submit an integrated summary of all available information about the safety of the drug product, including pertinent animal data, demonstrated or potential adverse effects of the drug, clinically significant drug/drug interactions, and other safety considerations, such as data from epidemiological studies of related drugs. The safety data shall be presented by gender, age, and racial subgroups. When appropriate, safety data from other subgroups of the population of patients treated also shall be presented, such as for patients with renal failure or patients with different levels of severity of the disease. A description of any statistical analyses performed in analyzing safety data should also be included, unless already included under paragraph (d)(5)(ii) of this section.*

The above referenced requirement of drug submissions suggests that an integrated analysis database derived from several studies that are included as part of the total submission package need to be created. Basically, the Integrated Summary of Safety (ISS) or Integrated Analysis of Safety (IAS) is in general an essential component of a submission. Within the Common Technical Document (CTD), the analysis results

derived from integrated analysis are located in Section 5.3.5.3. This information will then typically filter out to Sections 2.7.4 and 2.5 of the CTD. Table 1 summarizes the CTD placement and requirements of the information derived from integrated data.

**Table 1:** Placement and requirements of ISS information in the CDT

| | CTD Section | US Regulation | Comment |
|---|---|---|---|
| **2.5** | **Clinical Overview (~30 pages)** 2.5.4 Overview of Efficacy 2.5.5 Overview of Safety | N/A | Not a US requirement but recommended by ICH M4E |
| **2.7** | **Clinical Summary (~50 – 400 pages)** 2.7.3 Summary of Clinical Efficacy 2.7.4 Summary of Clinical Safety | 21 CFR 314.50(c)(2)(viii) | US requirement for a clinical summary |
| **5.3** | **Clinical Study Reports** 5.3.5.3 Reports of Analyses of Data from More than One Study (Including Any Formal Integrated Analyses, Meta Analyses, and Bridging Analyses) | 21 CFR 314.50(d)(2)(v) 21 CFR 314.50(d)(2)(vi) | Integrated Summary of Effectiveness Integrated Summary of Safety |

## 2. The Integrated Summary of Safety

The ISS differs from study level analysis due to the larger amount of data. The creation of integrated analysis database to support ISS work can be a challenging task as it requires harmonization and conversion of individual study data to the same format before initiation of programming and analysis. Analyses are usually performed on predefined groupings (pooling) of studies with common elements based on the integrated analysis safety analysis plan describing the pooling strategy. Differences between study designs, treatment, and duration of exposure and volume of data can create challenges in data pooling. Usually, the general idea is to combine data from all studies in some fashion and summarize the data as if they came from the same source. The summary will typically be based on some definition of treatment exposure and/or breakdown described in the integrated analysis statistical analysis plan (SAP). From an ISS reporting perspective, the results of all clinical studies performed on the investigational product are summarized as one single report based on the combined data.

There are several reasons why it is recommended to integrate data. This includes improving the precision of the incidence estimates, especially for rare events. Another reason is to enable the assessment of trends in small subgroups of patients, such as the elderly or special patient populations, where it may not be possible with study level data. Pooling results in a database that is much larger than the individual study databases presents a better chance of detecting potential safety problems that might have been missed in a small dataset. For example, one case of a safety concern here and there in a
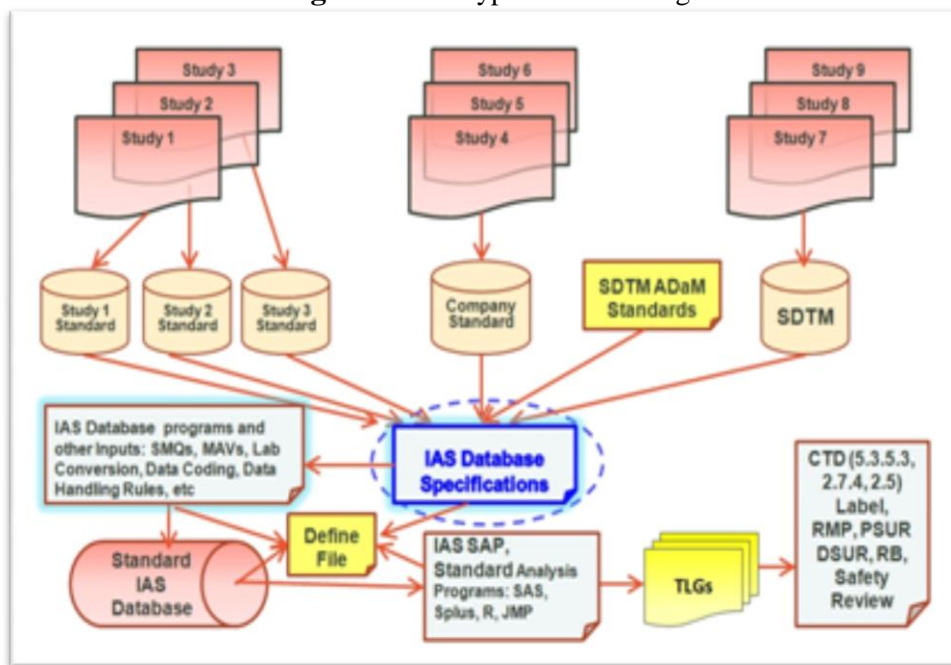
few individual studies will more clearly stand out when the data are combined together in a pooled database.

However, integrating data should be a well thought out process. Caution should be taken when interpreting integrated data as pooling data can lead to potential challenges in conclusions of cumulative information, for example regression to the mean. The merits of standardized integrated data for the purpose of conducting ISS work are immense. For example, it leads to improved efficiency in the preparation of submission documents. The adoption of a standard data format also allows for the opportunity to develop safety assessment tools that are reusable across compounds. With these tools, more time can be allocated to the team to focus on the interpretation of the results. This also can help in moving towards a strategic approach for improving the quality, speed, and transparency of submission documents.

## 3. The Typical Clinical Analysis Database Setting

The typical analysis data setting is one in which source clinical data are available from several studies, typically in a SAS format with some data potentially in non-SAS format. Further, study level data may be based on the Clinical Data interchange Standards Consortium (CDISC) Study Data Tabulation Model (SDTM) or they can be based on a company standard or some other non-SDTM standard or non-company standard. Figure 1 illustrates the typical data setting.

**Figure 1:** The typical data setting

The typical data domains within the SDTM are provided in Figure 2 below.

**Figure 2:** SDTM Data setting



Company data standards will typically follow this model or use similar data domain components with common elements to this model. For example, it is common to have demographics (DM) or adverse event (AE) datasets though the dataset names or variable names may not be the same as those used in SDTM.

In general from the point of view of creating an integrated analysis dataset, the first step is to create a standard set of comprehensive analysis specifications for the ISS database which are used to guide programming activities for creating the integrated analysis database. The integrated data specifications should accommodate the planned analyses discussed in the ISS SAP. They should also take into account any auxiliary input information, such as, identification of adverse events of special interest (AESIs) which are typically based on standardized Medical Dictionary for Regulatory Activities queries (MedDRA SMQs) or custom MedDRA queries.

The analysis data specifications should also incorporate information on markedly abnormal values (MAVs) or Potentially Clinically Concerning (PCCs) laboratory, vital signs, and electrocardiogram (ECG) data values. The MAV criteria can be custom criteria based on a company standard or they based on some other standard such as the Common Terminology Criteria for Adverse Events (CTCAE) Grades. The analysis database specifications must consider applying common standard derivation rules and data handling rules across studies for the purpose of creating a harmonized derivation rules.

## 4. The Standard Integrated Summary of Safety Database

### 4.1 Some Pros and Thoughts on a Standardized Integrated Analysis Database

There are many reasons why a standardized integrated analysis database should be developed. For instance, it can help to improve efficiency in the production of tables, listings, and graphs (TLGs). It can be used to facilitate for the development of common tools and utilities for analysis and also, it can help in promoting consistency and quality and can lead to ease of validation and the creation of analysis datasets that are submission ready. A standardized analysis dataset can also help to ensure program code quality, consistency, portability, reusability and also facilitate for easy maintenance. Another benefit of a standardized analysis database is that can lead to a reduction in programming

time and a decrease in the learning curve and an increase in productivity on routine analyses.

With the creation of standard integrated analysis database, more time can be allocated for data review which can permit for resource allocation and flexibility. If properly planned, a standard integrated analysis dataset can be streamlined and eliminate unnecessary variables from datasets. Careful thought should be given in the harmonization of data properties and variable attributes and ensure variable attribute consistency, that is, variable type, length, labels, and formats across datasets to ensure ease of use and future data additions.

## 4.2 Guiding Principles and Considerations for a Standard Analysis Database

A key principle of the analysis database is that it should be *analysis-ready*. However, the degree to which a dataset is considered to be analysis ready is not a clear cut notion or concept. In general, the right balance must be implemented to ensure that unnecessary programming time is not spent deriving analysis variables that can easily be created within the statistical analysis program and vice-versa, whether it is a descriptive or inferential analysis. A well designed integrated analysis database should also be useful for a variety of purposes, including integrated analysis of safety analyses, labelling, periodic safety reviews or safety updates, internal safety reviews and other regulatory needs. Further, the integrated analysis database must be built in such a way that additional studies can be added as needed. The integrated database must also have submission ready datasets and facilitate for easy documentation.

When one considers the points discussed above, an *ADaM-compliant* integrated analysis database provides the premises or reference through which such a database can be built. That is, in order to create a useful and viable integrated analysis database, a comprehensive set of data specifications for the ISS database based on the ADaM Basic Data Structure (BDS) is recommended. The specifications should focus on key or primary set of datasets and variables and eliminate unnecessary information. The ISS database specifications should account for technical considerations and rules for generating the analysis variables, for example, common algorithms for specific variables and data handling rules.

Careful thought must be given regarding the ease of development of the define files. Other key considerations include MedDRA and WHODrug versioning, laboratory data conversions, adverse events of special interest, and markedly abnormal values. Per the ADaM recommendation: *The overall principle in designing Statistical Analysis Datasets and related metadata is that there must be clear and unambiguous communication of the content, source and quality of the datasets supporting the statistical analyses performed in a clinical study. An analysis dataset serves as a central depository of raw data and analyzable variables derived from one or more of the raw datasets. The derived variables are used as inputs to produce statistical summaries*.

These principles should be taken into consideration when putting together an integrated analysis database specification. The CDISC recommendations on the use of the ADaM database principles, such as, redundancy for easy analysis in the sense that common variables may be found across all analysis datasets should apply. For example, population flags can be included in all the datasets. Further, the analysis datasets may have a greater number of numeric variables, for example, SAS formatted dates and numeric representation of a character grouping variable. The analysis datasets may also combine

variables from multiple data domains which may contain one or more records per subject, per analysis parameter and per analysis time point. Whatever the case, the specification must ensure that they fully describe the data being derived and include details of source variable that support the analysis datasets and include enough information to facilitate traceability.

## 5. Specific Considerations and Challenges in Creating Standard Integrated Analysis Database Specifications

Some specific considerations need to be addressed in creating the standard analysis database specifications. These include, but not limited to:

- Harmonization of dataset and variable attributes and ensure variable attribute consistency, that is, variable type, length, labels, and formats across datasets

- Identification of required/key variables across datasets

- Inclusion of flags in datasets – too much versus too little

- Imputation of data, such as missing date of last dose or changing character values to numeric, such as '< 23' to a numeric value

- MAV criteria for Labs, Vital Signs, and ECG parameters

- Traceability – maintain a solid relationship with source data with any discrepancies noted

- Keeping integration simple to accommodate future updates

- Comprehensive treatment coding

- General rules on titration studies and open label and *randomized-to-open-label studies*

- Determination of general guidelines for naming convention for subject level variables and derivations, for example, race categories

- Ensure transparency, reproducibility, accuracy, consistency, and ease of use

- Guidelines for flag variables, such as flags for adverse events of special interest (AESIs), concomitant medication classes, medical history classes

- Definitive MedDRA and WHODrug versioning processes

- Determination of cut-off dates for adverse events, laboratory data, vital signs, and ECG data.

## 6. Resolution on Some Specific Items

Various specific items of concern, including those noted above, need to be addressed accordingly in order to arrive at an acceptable resolution on how these will be addressed when creating the integrated analysis datasets. These resolutions must be detailed in the specifications or elsewhere as appropriate. Some of the possible resolutions to some of these considerations include:

- Only those flags that are needed across more than one table, for example, population flags should be included in each dataset

- Use the principle of *ADaM on a diet*, that is, eliminate unnecessary flags or flags that are only used once in an analysis or those that can be easily derived within the analysis program

- If flags are used across different datasets, they should only be derived once in one dataset and *forward-populated*

- Flag variable names and labels should be named appropriately and meaningfully, for example, one can use TEAEFL instead of FLAG1

- Similarly for flags for special interest adverse events, use descriptive variable names, such as, AEDescripFL where Descrip is a meaningful descriptor of the cluster. For example, use: AEMACEFL for MACE events; CMDescripFL where Descrip is a meaningful descriptor of a concomitant medication class, for example, CMDIABFL for Diabetes Medications; MHDescripFL where Descrip is meaningful descriptor of the medical history, for example, MHDIABFL for History of Diabetes

- In the SAS dates, the time component can be left out unless it is used to determine occurrence events or data collection relative to dosing time

- Any data imputation should include a flag indicating that the specific value is imputed

- MAV criteria should be based on standard set of criteria.

- All analysis data derivations descriptions should be specific enough to allow for traceability

- Use standard rules for imputing missing date of first dose and last dose

- For titration studies, dosing information should include details of the dosing dates for each segment of dose

- Subject level variables should retain the same names from the source datasets where there are derived from

- Further clarify variable names as necessary and label names as appropriate, for example, SBPBL can be used to clarify that this is Systolic Blood Pressure at Baseline as opposed to just using: SBP – Systolic Blood Pressure

- The development of AESIs clusters should use SMQs or custom based versions that are well documented

- The cut-offs dates for AEs should be standard, but can be adjusted as needed, for example, drugs with long half life

- Imputation of missing or partial dates should follow specific data handling rules

- For lab data, vital signs data, and ECG, cut-off should always be fixed, for example, 7 days after last dose, but can be changed as appropriate, for example, for drugs with long half-life

- When applying flags in the process of converting lab data from one unit to another, carry the flag from the value that was originally assigned to the original lab data.

## 7. Other Considerations

From then point of view of the ISS, the usual key or core safety datasets need be considered in the analysis database specifications. These will typically include:

- **ADSL** - Subject Information
- **ADSD** - Study Drug
- **ADMH** - Medical History
- **ADCM** - Concomitant Medications
- **ADAE** - Adverse Events
- **ADLB** - Clinical Laboratory
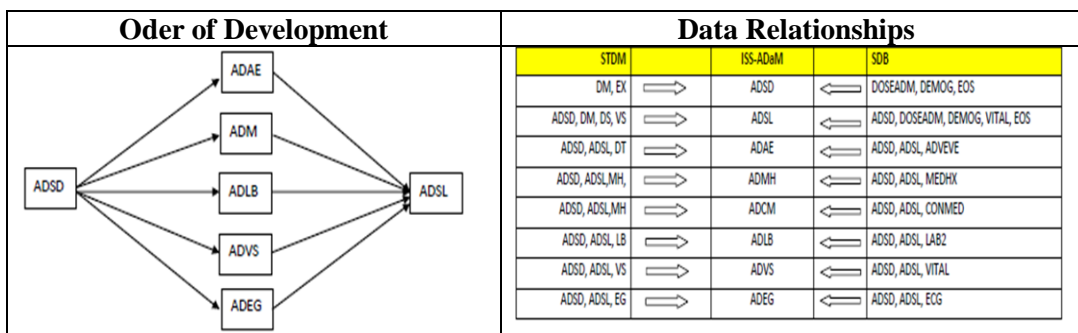- **ADEG** – Electrocardiogram
- **ADVS** – Vital signs

From a programming perspective, it is ideal that the integrated analysis database specification address sources of data that are both SDTM compliant and company standard based if available or a mixture of the different source formats. The specifications must clearly allow for traceability in all cases. One possibility is to build analysis specifications that describe how variables are derived using difference sources. A possible way in which this can be done is given in the snapshot in Table 2.

**Table 2:** Snapshot of a Possible Data Specification

| ADSL (Subject Information) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Note: Include all subjects randomized, including subjects not dosed | | | | | | | | |
| | | | | | SDTM | | SDB | |
| Variable Name | Label | Type | Controlled Terms/ Format | Core | Source | Derivation Algorithm/ Comments | Source | Derivation Algorithm/ Comments |
| STGRPN | Study Group (N) | Num | STGRPF | Perm | Derived | Derived based on analysis pools | Derived | Derived based on analysis pools |
| STGRP | Study Group | Char | | Perm | Variable STGRPN | Derived from STGRPN | Variable STGRPN | Derived from STGRPN |
| STUDYID | Study Identifier | Char | | Req | From DM, variable STUDYID | | From DEMOG, variable STUDY | |
| USUBJID | Unique Subject Identifier | Char | | Req | From DM, variable USUBJID | Must be unique for a subject throughout a submission and be derived based on the first study a subject was enrolled in. | From DEMOG, variables STUDY, SITE, and SUBJNO | Must be unique for a subject throughout a submission and be derived based on the first study a subject was enrolled in. Equals compress(STUDY‖'-'‖put(SITE,best.)‖'-'‖put(SUBJNO,best.)) |
| SUBJID | Subject Identifier for the Study | Char | | Req | From DM, variable SUBJID | | From DEMOG, variable SUBJNO | Character representation of SUBJNO, including leading zeros |
| SUBJNO | Subject Identifier for the Study (N) | Num | | Req | From DM, variable SUBJNO | | From DEMOG, variable SUBJNO | |
| SITEID | Study Site Identifier | Char | | Req | From DM, variable SITEID | | From DEMOG, variable SITE | Charater representation of SITE, including leading zeros |
| AGE | Age | Num | | Req | From DM, variable AGE | If age is missing, calculate from informed consent | From DEMOG, variable AGE | If age is missing, calculate from informed consent date |
| AGEU | Age Units | Char | AGEU | Req | From DM, variable AGEU | | From DEMOG, variable AGEUNIT | Derived from AGEUNIT |
| BRTHDT | Date of Birth | Num | DATE9. | Req | From DM, variable BRTHDTC | Date substring of BRTHDTC | From DEMOG, variable BIRTHDT | |

In terms of order of development, as noted earlier, it recommended that the forward population principle be used. One approach is to develop the ADSD dataset first and the others thereafter as the other datasets will typically use data from ADSD. Figure 3 below show how this can be done and how the order of development may proceed and how the source datasets feed into domains of the integrated analysis database.

**Figure 3**: Order of development and source data and analysis data relationships



| Oder of Development | Data Relationships |
|---|---|

| STDM | | ISS-ADaM | | SDB |
|---|---|---|---|---|
| DM, EX | ⟹ | ADSD | ⟸ | DOSEADM, DEMOG, EOS |
| ADSD, DM, DS, VS | ⟹ | ADSL | ⟸ | ADSD, DOSEADM, DEMOG, VITAL, EOS |
| ADSD, ADSL, DT | ⟹ | ADAE | ⟸ | ADSD, ADSL, ADVEVE |
| ADSD, ADSL,MH, | ⟹ | ADMH | ⟸ | ADSD, ADSL, MEDHX |
| ADSD, ADSL,MH | ⟹ | ADCM | ⟸ | ADSD, ADSL, CONMED |
| ADSD, ADSL, LB | ⟹ | ADLB | ⟸ | ADSD, ADSL, LAB2 |
| ADSD, ADSL, VS | ⟹ | ADVS | ⟸ | ADSD, ADSL, VITAL |
| ADSD, ADSL, EG | ⟹ | ADEG | ⟸ | ADSD, ADSL, ECG |

**Note:** The datasets with non-SDTM names represent possible names for data domains of company standard.

## 8. Conclusion

In this paper, we attempted to describe an approach that can be taken to create a harmonized integrated analysis database from both SDTM and non-SDTM compliant data sources based on ADaM-compliant BDS analysis model. The primary idea to do this through the development of a comprehensive set of standard analysis database specifications that are used to drive programming activties. Such an effort can lead to the

creation of integrated analysis databases that can help in ensuring coding efficiency quality, consistency, portability, reusability and easy maintenance. This can also help to facilitate for the development of common tools and utilities for analysis thereby improving efficiency and quality and making it easy to validate outputs. Additionally this can help promote cross-study and cross-compound analysis and standardization and in optimizing data processing and reporting activities across submissions. It can also lead to savings in programming time and increase productivity.

# References

ABT, K. and Krupp, P. (1986),"Pooling of laboratory data safety data in multicenter studies,", *Drug Information Journal*, **20**, 311–313.

Barrows, C. (2011),"Principles and practicalities in building ADaM datasets," *PhUSE Single Day Events*.

Case, T. (2011), "Adverse events of special interest and MedDRA Upgrades: A dilemma and proposed solution," *PharmaSUG20XX*, Paper TT09.

CDISC - "Analysis Data Model Implementation Guide".

CDISC, "Analysis Data Model (ADaM)," *CDISC The ADaM Data Structure for Adverse Event Analysis*.

Chang, S. and Wong, S. (2010), "Successful FDA advisory meetings using analysis datasets," *PharmaSUG2010*.

Chen, H-L, and Wang, H. (2012), "Multiple application of ADaM time-to-event datasets," *PharmaSUG 2012*, Paper DS19.

Cui, X., Chen, M., and Moseley, M. (2011), "New tips and tricks for creating a harmonized, report-friendly SDTM and ADaM lab data for a clinical study report," *PharmaSUG2011*, Paper TU02MA.

Decker, C., Zhao, Y., and Hirschfeld, S. (2010), "The CDIS/FDA integrated data pilot: a case study in implementing CDISC standards to support an integrated review'" *PharmaSUG* , Paper RS02.

Dirk, V. K (2011), "ADaM on a diet: preventing wide and heavy analysis datasets," *PhUSE 2011*, Paper CD11.

Enas, G. and Goldstein, D. (1999), "Defining, Monitoring and Combining safety information in clinical trials," *Statist. Medic.*, **14**, 1099–1111.

Endri, R. H. (2011), "Much ADaM about nothing a proc away in a day," *Pharmaceutical Programming*, **4**, 45–50.

FDA: Guidance for The Format and Content of the Clinical and Statistical Sections of New Drug Applications.

*Section H: Integrated Summary of Safety Information*.

FDA Guidance for Industry (2009), "Integrated Summaries of Effectiveness and Safety: Location Within the Common Technical Document".

Fissekis, J. A. (2009), "Integrated summaries of safety and efficacy", In L. F. Wood and M. Foote (eds), *Targeted Regulatory Techniques: Clinical Documents for Drugs and Biologics*, **Chapter 9**, 1250–130.

Fulton, J. (2010), Dealing with lab data stacking the deck in your favor. *SAS Global Forum*, Paper 182–2010.

Huang, J. H. (2010), "How to build ADaM from SDTM - a real case study", *PharmaSUG2010*, Paper CD06.

Guidance for Industry (2009) "Integrated summaries of effectiveness and safety - Location within the common technical document".

Kenny, S. J. and Helton, E. D. (2006), "Using CDISC models for the analysis of safety data", *Phuse2006*, Paper RA02.

Kenny, S. J. (2011), "Using CDISC models for the analysis of safety data," *PhUSE2006*, Paper RA02.

Kenny, S. J. (2008), "Integrated Databaset," *In Encyclopedia of Clinical Trials*.

Kenny, S. J. and Litzsinger, M.A. (2010), "Strategies for implementing SDTM and ADaM standards," *PharmaSUG*, Paper FC03.

Kubler, L. and Weihrauch, T. R. (2002), "Integrated summaries and meta-analyses in clinical drug development," *Drug Information Journal*, 36, 127–133.

Li, D., Li, S., and Sproule, S. (20XX), "A SAS macro for creating adverse event analysis dataset," *PharmaSUG20XX*, PaperCC06.

Lievre, M. Cucherat, M., and Leizorovicz (2002), "Pooling, meta-analysis, and the evaluation of drug safety", *Current Controlled Trials in Cardiovascular Medicine*, **3**, 1–4.

Lin, E. Z. and Ding, B. (2012)' Analysis-ready considerations, implementation, and real world applications,"

*PharmaSUG 2012*, Paper DS13.

Matthews, C. (2009), "Techniques for assigning NCI CTC Grades to laboratory results," *PhUSE 2009*, Paper PO07.

Matthews, C. (2009), "Assigning NCI CTC Grades to laboratory results," *PharmSUG2010*, Paper PO03.

McEntegart, D. J. (2000), "Pooling in integrated databases," *Drug Information Journal*, **34**, 495–499.

Meyerson, L. J. (2003), "Integrated summary report," *In Encyclopedia of Biopharmaceutical Statistics*, Pages 486–489.

Peterson, T. and Izard, D. (2010), "The 5 biggest challenges of ADaM," *PharmaSUG2010*, Paper CD10.

Ryan, P. J. D. (1995), "Integrated clinical databases: Developing and effective strategy towards CANDA,"*Drug Information Journal*, **29**, 767–771.

Saranadasa, P. (2009), "Successful lab result conversion for LAB analysis data with minimum effort," *NESUG*.

Saranadasa, P. (2011), "Challenges in implementing ADaM datasets: Balancing the analysis-ready and traceability concepts," *PharmacSUG2011*, Paper CD19.

Sharma, R. (2012), "Creating an integrated summary database using CDISC ADaM: Challenges and tips, and things to watch out," *PharmaSUG 2012*, Paper DS17.

Smith, D.J., Schulz, D., Kloss, G., and Cheng, W. (2010), "Considerations for building an integrated safety database using SAS," *PharmaSUG2010*, Paper AD15.

Steffensen, K. F. and Jepsen, G. D. (2009), "An implementation of ADaM standards not driven by a submission," *PhUSE 2009*, Paper CD13.

Temple, R. J. (1991), "The regulatory evolution of the integrated safety summary," *Drug Information Journal*, **25**, 485–492.

Temple, R. (2006), "CTD ISS/ISE introduction and Summary of Issues".

Van Bemmelen, J. (2008), "Applying SMQs to adverse events data," *PhUSE2008*, Paper TU05.

Wang, Q. and Herremans, C. (2010), "How to go from an SDTM finding domain to an ADaM-compliant Basic Data."