


Biostatistics, Epidemiology &
Research Design (BERD)

Howard Cabral, PhD, MPH
Christine Chaisson, MPH

Seminar Series: CTSI Presents...


December 15, 2016



CTSI Presents...

**Data Management:
What you should know
and why you should care**


Christine Chaisson
Director, Data Coordinating Center
Assistant Research Professor, Biostatistics
Boston University School of Public Health



Boston University School of Public Health
Data Coordinating Center


Things that can go wrong with data

- Crucial data elements may be missing
- Data may be incorrect due to errors in:
 - Data collection
 - Data entry
- Data may not have common identifier
 - Cannot be merged
 - May be merged incorrectly
- Data may not be saved or backed up
- Data files may be lost or corrupted



Real World Examples

- A few illustrations of common data problems from the popular news sources



The Inquirer DAILY NEWS

philly.com August 3, 2012


Forbes: Bad data hurt Haverford in college rankings

"Forbes' annual list is out, and Haverford plummeted from No. 7 to No. 27 - for no obvious reason. A College spokesman explained that the error was based on single figure:

A zero was incorrectly entered in database instead of 108 for the graduation rate of white women who enrolled in 2004.

...But no revision is planned, since the magazine and the online list has already been published."

Data Entry Error



PharmaTimes ONLINE May 6, 2012

Vertex stock slides over cystic fibrosis data mistake

"Shares in Vertex Pharmaceuticals have taken a hit after the company had to take the rather embarrassing step of correcting previously-announced interim mid-stage results of a combination cystic fibrosis treatment.

...the result of a misinterpretation of the denominator of the treatment group between the firm and its outside statistical ve

Data Mismanaged



Gene: Excel Error Calls Into Question... Posted 22 Apr 2013

SPECTRUM

- 'This Time It's Different' a 2009 book by Harvard researchers (Reinhard & Rogoff) contained "Serious errors that inaccurately represent the relationship between public debt and GDP growth among 20 advanced economies in the post-war period." Investigators at UMass (among others) were unable to replicate work which led to **accusations of intentional fudging of the data or**
- The Authors admitted they failed to include five rows of data from an Excel file (Australia, Austria, Belgium, Canada, and Denmark) —a "coding error" which they said was "a significant lapse on our part"

excluded key data

The New York Times July 7, 2016
How Bright Promise in Cancer Testing Fell Apart

- Duke Cancer Center's gene-based tests proved worthless, research behind them was discredited
- Statisticians from MD Anderson discovered errors such as columns moved over in a spread-sheet; Duke team "shrugged them off" as "clerical errors."
- Four papers were retracted
- Duke shut down three cancer trials
- Center leaders resigned or were removed
- People died and their relatives sued Duke

Data Entry Management in Excel

Goal: Convert Data into Electronic Format as Quickly as Possible

How should data be managed?

- No single "right" way to collect or manage data
- Consider:
 - Environment/location
 - Resources (\$)
 - Regulations
- Be sure to **plan** prior to study start
- Do what works for the study at hand

Where to start?

Data management plan

From Wikipedia, the free encyclopedia

A **data management plan** or **DMP** is a formal document that outlines how you will handle your data both during your research, and after the project is completed.^[1] The goal of a data management plan is to consider the many aspects of data management, metadata generation, data preservation, and analysis before the project begins; this ensures that data are well-managed in the present, and prepared for preservation in the future.

DMP Purpose: To help you collect, manage and share your data; meet funder requirements. General elements include:

- Project or study description
- Documentation, organization, storage
- Access and sharing
- Archiving

BOSTON UNIVERSITY

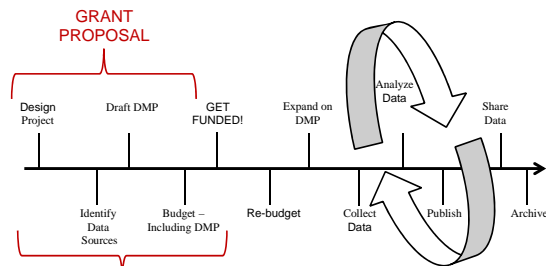


DMP: Basic Elements

- Study design, data types and sources
- Storage format and location
- Naming conventions, documentation
- Software used for manipulation
- Project Staff – who has permission to what
- Identifiers (if applicable)
- Back ups, security
- Archiving
- Sharing

BOSTON UNIVERSITY

Timeline



BOSTON UNIVERSITY

Beginning: Identify Key Data Elements

- Review hypotheses
- What are primary, secondary outcomes?
- What covariates and confounders must be collected?
- What are the data sources?
 - Questionnaires
 - Labs, imaging
 - Medical record review
 - other external sources (e.g., lab results, medical records, death certificates)

BOSTON UNIVERSITY

Other Data Elements

- Regulatory data:
 - IRB requirements
 - Safety (DSMB)
 - FDA (e.g., 21 CFR, part 11)
 - VA
 - Other?
- Tracking/Study management data:
 - Tracking participants
 - Tracking data elements by time-points
- Harmonization and sharing
 - NIH
 - Other

BOSTON UNIVERSITY

Visit Protocol: Data by Time-point

- Determine visit Schedule and data collected at each visit
 - Questionnaires
 - Labs
 - Other?
- Consider data not be connected to visits
 - Adverse events, serious adverse events
 - Hospitalization
 - Death
 - Medical records



BOSTON UNIVERSITY

Sample data/visit grid

Screening Period	Baseline	Week 4	Week 8	Week 12	Week 16	Week 20	Week 24	Follow-up visit week 26
1:1 visit between Day 28 to 31	Day 1	Day 28	Day 56	Day 84	Day 112	Day 140	Day 168	Day 200 visit
Medical history & demographic (Form 100)	X							
Physical examination (PE) (1)	X							
Physical examination (PE) (2)		X	X	X	X	X	X	X
Emergency visit / Transfer (1)	X	X	X	X	X	X	X	X
Inclusion / exclusion criteria evaluation	X							
Informed consent / assent	X							
Randomization	X							
Targeted physical exam (1)		X	X	X	X	X	X	X
Targeted physical exam (2)			X	X	X	X	X	X
Site visit with joint sponsor (S1) (1)								
Site visit with joint sponsor (S1) (2)								
Site visit with joint sponsor (S1) (3)								
Site visit with joint sponsor (S1) (4)								
Site visit with joint sponsor (S1) (5)								
Site visit with joint sponsor (S1) (6)								
Site visit with joint sponsor (S1) (7)								
Site visit with joint sponsor (S1) (8)								
Site visit with joint sponsor (S1) (9)								
Site visit with joint sponsor (S1) (10)								
Site visit with joint sponsor (S1) (11)								
Site visit with joint sponsor (S1) (12)								
Site visit with joint sponsor (S1) (13)								
Site visit with joint sponsor (S1) (14)								
Site visit with joint sponsor (S1) (15)								
Site visit with joint sponsor (S1) (16)								
Site visit with joint sponsor (S1) (17)								
Site visit with joint sponsor (S1) (18)								
Site visit with joint sponsor (S1) (19)								
Site visit with joint sponsor (S1) (20)								
Site visit with joint sponsor (S1) (21)								
Site visit with joint sponsor (S1) (22)								
Site visit with joint sponsor (S1) (23)								
Site visit with joint sponsor (S1) (24)								
Site visit with joint sponsor (S1) (25)								
Site visit with joint sponsor (S1) (26)								
Site visit with joint sponsor (S1) (27)								
Site visit with joint sponsor (S1) (28)								
Site visit with joint sponsor (S1) (29)								
Site visit with joint sponsor (S1) (30)								
Site visit with joint sponsor (S1) (31)								
Site visit with joint sponsor (S1) (32)								
Site visit with joint sponsor (S1) (33)								
Site visit with joint sponsor (S1) (34)								
Site visit with joint sponsor (S1) (35)								
Site visit with joint sponsor (S1) (36)								
Site visit with joint sponsor (S1) (37)								
Site visit with joint sponsor (S1) (38)								
Site visit with joint sponsor (S1) (39)								
Site visit with joint sponsor (S1) (40)								
Site visit with joint sponsor (S1) (41)								
Site visit with joint sponsor (S1) (42)								
Site visit with joint sponsor (S1) (43)								
Site visit with joint sponsor (S1) (44)								
Site visit with joint sponsor (S1) (45)								
Site visit with joint sponsor (S1) (46)								
Site visit with joint sponsor (S1) (47)								
Site visit with joint sponsor (S1) (48)								
Site visit with joint sponsor (S1) (49)								
Site visit with joint sponsor (S1) (50)								
Site visit with joint sponsor (S1) (51)								
Site visit with joint sponsor (S1) (52)								
Site visit with joint sponsor (S1) (53)								
Site visit with joint sponsor (S1) (54)								
Site visit with joint sponsor (S1) (55)								
Site visit with joint sponsor (S1) (56)								
Site visit with joint sponsor (S1) (57)								
Site visit with joint sponsor (S1) (58)								
Site visit with joint sponsor (S1) (59)								
Site visit with joint sponsor (S1) (60)								

- ### Timelines and Tasks
- Develop Protocol and Analytic Plan
 - Create and pilot of forms/assessments
 - Design/construct data systems
 - Data Collection/entry
 - Participant/Data Tracking
 - Subject recruitment
 - Data collection (baseline and follow-up)
 - Data cleaning, auditing, and QA
 - Preliminary analysis
 - Manuscript preparation & submission
-

- ### Create a Visual Timeline
- It doesn't have to be fancy
 - More detail is better but something simple is better than nothing
 - Plan to review and revise it often
-

Simple overview Timeline

Study Timeline

Selected Activities	1-6	12-18	24-30	36-42	48-54
Hiring & training	X				
Finalize instruments & IRB	X				
Enrollment	X	X	X	X	
Intervention		X	X	X	X
Follow-up		X	X	X	X
Data QA/clean		X	X	X	X
Primary & secondary data analyses					X
Presentations & Publication					X
Study meetings	X	X	X	X	X

Sample Task-based Gantt

Tasks	Year 1				Year 2			
	Months 1-3	Months 4-6	Months 7-9	Months 10-12	Months 13-15	Months 16-18	Months 19-21	Months 22-24
Finalize CRFs	█							
IRB Approval	█							
Finalize data platforms	█							
Finalize protocol	█							
Build eCRF	█							
Build database	█							
Pilot CRFs/protocol		█						
Build website		█						
Enrollment/data collection			█					
Query data / monitor				█				
Automate data reports					█			
Update website, reports						█		
OSIRIS data freeze, reports, meeting							█	
Follow up visits								█

Multi-task Indicating Responsible Parties

Task	Responsible Party	Start	End
Finalize CRFs	DP	Dec 12	Jan 12
IRB Approval	DP	Dec 12	Jan 12
Finalize data platforms	DP	Dec 12	Jan 12
Finalize protocol	DP	Dec 12	Jan 12
Build eCRF	DP	Dec 12	Jan 12
Build database	DP	Dec 12	Jan 12
Pilot CRFs/protocol	DP	Dec 12	Jan 12
Build website	DP	Dec 12	Jan 12
Enrollment/data collection	DP	Dec 12	Jan 12
Query data / monitor	DP	Dec 12	Jan 12
Automate data reports	DP	Dec 12	Jan 12
Update website, reports	DP	Dec 12	Jan 12
OSIRIS data freeze, reports, meeting	DP	Dec 12	Jan 12
Follow up visits	DP	Dec 12	Jan 12

Tools Of The Trade

- Analytic plan
- Detailed protocol
- Well designed data collection forms
- Tracking system
- Data capture/entry system
- Plan for data query (checking/cleaning)
- Manuals
- Data dictionaries



Web Site User Manual
Version date: August 28, 2010
<http://www.tbkweb.bumc.bu.edu/SIZANANI/>

TABLE OF CONTENTS

- I. Introduction 1
- II. Purpose of the Website 1
- III. Data Management 1
- IV. Site Use of the Website 1
- V. Security 2
- VI. Website Navigation 3
- VII. Data Entry 3
 - A. Selecting Study Site 4
 - B. Data Entry Guidelines 4
 - 1. ID Number 4
 - 2. Order of Forms 4
 - 3. Case Fields 4
 - 4. Time Fields 4
 - 5. Consent Fields 4
 - 6. Blank Fields 4
 - 7. Participation of Random Buttons 4
 - 8. Copying the DOC 4
 - C. Data Field Types 4
 - D. Skip Patterns 4
 - E. Error Messages 4
 - F. Saving Data 4
- VIII. Logging on to the SIZANANI Website 5
- IX. Main Menu 10
 - A. Data Entry 11
 - 1. Screening Log 11
 - 2. Screening Log Follow-up 11
 - 3. ICU Contact Information 11
 - 4. ICU Navigator Contact Log 11
 - 5. Hospital Discharge Form 11
 - 6. HIV Knowledge Resources 11
 - 7. TB Knowledge Resources 11
 - 8. ICL Participation Log 11
 - 9. ICL Knowledge Update Log 11
 - 10. ICL Study Conclusion 11
 - B. Study Information and Documents 11
 - C. Reports 11
 - 1. Enrollment Report 11
 - 2. Study Data Summary 11
 - E. Training 11
 - F. Follow-Up Visit List 11
- X. Contact Information 29



Consider Languages

U19 STUDY
U19 研究

MANUAL OF PROCEDURES
程序手册

Version date: March 20, 2011
版本日期: 2011年3月20号

中国医学科学院皮肤病研究所
中国疾病预防控制中心性病控制中心
NATIONAL CENTER FOR STD CONTROL, CHINA CDC

TABLE OF CONTENTS (目录)

- I. INTRODUCTION 简介 4
- II. CONTACT INFORMATION 联系方式 4
- III. STUDY PROCEDURES 研究程序 4
- A. U19 U19 研究程序 4
- B. U19 U19 研究程序 4
- C. U19 U19 研究程序 4
- D. U19 U19 研究程序 4
- E. U19 U19 研究程序 4
- F. U19 U19 研究程序 4
- G. U19 U19 研究程序 4
- H. U19 U19 研究程序 4
- I. U19 U19 研究程序 4
- J. U19 U19 研究程序 4
- K. U19 U19 研究程序 4
- L. U19 U19 研究程序 4
- M. U19 U19 研究程序 4
- N. U19 U19 研究程序 4
- O. U19 U19 研究程序 4
- P. U19 U19 研究程序 4
- Q. U19 U19 研究程序 4
- R. U19 U19 研究程序 4
- S. U19 U19 研究程序 4
- T. U19 U19 研究程序 4
- U. U19 U19 研究程序 4
- V. U19 U19 研究程序 4
- W. U19 U19 研究程序 4
- X. U19 U19 研究程序 4
- Y. U19 U19 研究程序 4
- Z. U19 U19 研究程序 4

Sample Data Dictionary: SAS formatted dataset

Number	Variable	Label	Type	Format
1	idptg_mch_id	ADEPT Study ID	Num	8
2	idptg_mch_id	Study Assessment	Num	11
3	idptg_mch_id	Study Site	Char	150
4	idptg_mch_id	Parent Language	Char	520
5	idptg_mch_id	Completed	Num	11
6	idptg_mch_id	Discontinued	Num	DATEMMYY3
7	idptg_mch_id	Discontinued	Num	DATEMMYY3
8	idptg_mch_id	When is study start?	Num	DATEMMYY3
9	idptg_mch_id	Please refer to the research assistant's initials	Char	\$101.
10	idptg_mch_id	Language of Interview	Num	LANSBP
11	idptg_mch_id	Gender	Num	SEXF
12	idptg_mch_id	Weight	Num	WTG
13	idptg_mch_id	Participate weight in kg	Num	WTG
14	idptg_mch_id	Age	Num	8
15	idptg_mch_id	DOB Day	Num	8
16	idptg_mch_id	DOB Month	Num	8
17	idptg_mch_id	DOB Year	Num	8
18	idptg_mch_id	Tube Site	Num	NOVEX
19	idptg_mch_id	Tube Site	Num	NOVEX
20	idptg_mch_id	Tube Site	Num	NOVEX
21	idptg_mch_id	Tube Site	Num	NOVEX
22	idptg_mch_id	Tube Site	Num	NOVEX
23	idptg_mch_id	Tube Site	Num	NOVEX
24	idptg_mch_id	Tube Site	Num	NOVEX
25	idptg_mch_id	Tube Site	Num	NOVEX
26	idptg_mch_id	Tube Site	Num	NOVEX
27	idptg_mch_id	Tube Site	Num	NOVEX
28	idptg_mch_id	Tube Site	Num	NOVEX



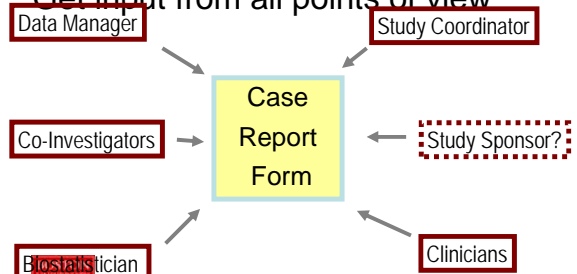
Data Collection Forms

- Data will usually end up on some type of "form" whether it is interview, chart review, or imaging results
- Make sure you plan carefully and leave time as this can be a lengthy process



Test Name	Result	Unit	Reference Range	Abnormal
RDW	13.7	g/L	11.5-14.5	***
WBC	2.8	x10 ⁹ /L	4.00-11.00	***
PLT	98	x10 ⁹ /L	140-400	***
MCV	81.1	fL	80.0-100.0	***
Neut	2.18	x10 ⁹ /L	2.00-7.00	***
Lymph	0.42	x10 ⁹ /L	1.00-3.00	***
Mon	0.2	x10 ⁹ /L	0.25-1.00	***
Baso	0.13	x10 ⁹ /L	0.01-0.10	***
Eos	0.19	x10 ⁹ /L	0.02-0.50	***
Rbc	0.19	x10 ¹² /L	4.10-5.10	***
Hct	0.110	L	0.380-0.480	***

Creating a Data Collection Form: Get input from all points of view



Real-world example why you should get input from multiple perspectives

What is your relationship to [person with illness] (i.e., [person with illness] is your _____)?

- Parent
- Child
- Sibling
- Spouse/Partner
- Other _____



Why is this topic important?

- Sloppy forms indicate sloppy research
- CRF may not answer study questions
- Danger of collecting:
 - too much data
 - too little data
 - the wrong data
- Annoyed:
 - Participants
 - Study Coordinator
 - Data Analyst...



When designing forms...

Think Google

Not Yahoo



What makes a good case report form?

- User-friendly, uncluttered, well organized
- Provides clear instructions for completion
- Terminology familiar to person filling out
- Reading level matches study participants/evaluators
- Unambiguous questions
- Questions only asked/data collected in one place and *only* one place
- Easy to refer back and clean data



Account For Missing Data

CBC	Unit	Value	
1. Hemoglobin	g/dl	___ . ___	<input type="checkbox"/> Not Done
2. Hematocrit	%	___ . ___	<input type="checkbox"/> Not Done
3. RBC	M/mm ³	___ . ___	<input type="checkbox"/> Not Done



Specify the Units

Alcoholic Beverages
Serving Sizes

A 12-ounce bottle or can of beer is 1 serving.
A 12-ounce glass of beer is 1 serving.
A 1-1.5 ounce shot of liquor straight or in a mixed drink is 1 serving.
A 5-ounce glass of wine is 1 serving.

Keep in mind that alcoholic drinks may contain more than one serving of alcohol.

How many servings of alcoholic beverages did you drink?

	1-24 Hours Preceding Gout Attack	25-48 Hours Preceding Gout Attack
*Beer	[Please Select]	[Please Select]
*Wine	[Please Select]	[Please Select]
*Spirits	[Please Select]	[Please Select]

Submit

Categorize Anticipated Responses

- ₁ USA
- ₂ Guatemala
- ₃ Mexico
- ₄ Dominican Republic
- ₅ Other
- USA
- Guatemala
- Mexico
- Dominican Republic
- El Salvador
- Other



ID Assignment

- Must be UNIQUE for each subject
- Should appear on every form (preferably page)
 - Primary key for records database
 - “Merge key” to join various data elements
 - Serves as identifier (no names!)
- May be a simple number 1001
- May be multi-part: 102101
 - 1 = Site
 - 02 = Language
 - 101= ID



Example of an ID that is not unique

The screenshot shows a form for 'Overdose Prevention & Naloxone'. A red circle highlights the 'Date' field containing the number '33015'. A blue box with the text 'Do NOT do this' is placed over the date field. Below the form, there are codes and abbreviations:

Codes and Abbreviations:
 FM Female to Male transgender
 MF Male to Female transgender
 NEP Code First three letters of mother's first name + date of birth (mm/dd/yy) Ex: GER053077
 BSAS Code First & third letters of first and last name Ex: Joseph "Joe" Francis Blow=JSBO

Don't Underestimate Need for Version# /Date

The screenshot shows a section of a form with two tables. The first table is for 'FOOT POSITION' and the second is for 'LEG LENGTH DISCREPANCY'. At the bottom of the form, there is a line of text: 'Partial meniscectomy vs. nonoperative management in meniscal tear with C... Version 1.0 & 01/05'. The 'Version' and 'Date' parts of this text are circled in red.

Consider who is Completing a Form

The screenshot shows a multi-step form titled 'Participant Forms'. The steps are: 'Scleroderma Health QOL', 'Coordinator Forms', 'Doppler Echocardiogram', and 'Renovascular Assessment'. A box labeled 'Clinician' is placed over the 'Doppler Echocardiogram' section, indicating that this section is completed by a clinician.

In Summary, when designing questions:

- Avoid ambiguous questions and open ended responses
- Include clear instructions
- Be sure form complexity matches collector (self, study coordinator, clinician)
- Collect data elements in correct format (“continuous” or “categorical”)
- Make categories mutually exclusive
- Pilot your forms in the target population



The Good The Bad And The Ugly

Paper or Electronic?

Paper Forms / Manual Entry

Advantages

- The old “standard”
- Shorter start-up time (Word/PDF)
- Relatively easy to train staff
- Hardcopy document to refer back to
- Can be done anywhere

Disadvantages

- Costs: data entry, storage and shipping
- Longer time from collection to database
- Errors in data collection (missing, out of range, skips)

Electronic Data Capture

Advantages

- Cleaner data at entry (required fields, skips, ranges)
- Can use data in real time (or close to it)
- No extra data entry costs
- Data can inform next visit even for short follow up

Disadvantages

- Programming time and costs
- Increased hardware and software costs
- Infrastructure concerns (software versions, internet connection, back-up equipment)
- Data security

A Word About “Canned” Software

- Many “canned” software packages available
- No single best choice
- Cost can vary widely
- Database structures vary
- Do your homework to make sure what for your project

Find out what software is available through your institution

Now I'm going to ask you about alcohol you took in the last 2 months.
In the past 2 months, how often did you have a drink containing alcohol?

Never
Monthly or less
2 to 4 times a month
2 to 3 times a week
4 or more times a week
Don't Know
Exclude

Назад Вперед

У нас есть несколько вопросов о вашем состоянии здоровья. Если вы сейчас живете? (пожалуйста, нажмите все подпадающие варианты)

Владели ли вы автомобилем? Или у вас есть мотоцикл? Или вы арендовали автомобиль? Или вы арендовали мотоцикл? Или вы арендовали мотоцикл? Или вы арендовали мотоцикл? Или вы арендовали мотоцикл?

На улице (пешеходы, забросанные камни, парки, и т.д.) Другие

Откажитесь от ответа.

Другие (укажите)

Consider languages when selecting data entry methods and software

What to use...?


- To determine what software is best suited for your project see:
 - What is available to you?
 - What is the cost (can you afford it)?
 - What has the features you need (e.g., language)

BOSTON UNIVERSITY

50

Once data are collected...


- Get your data into a useful format
- No "right" format – use what works for you
 - SQL database
 - SAS datasets
 - SPSS
 - Excel (be careful!)



BOSTON UNIVERSITY

Keep Personal Identifiers Separate

- Do not store any identifier unless you have a good reason for it
- Do not store identifiers in same files with study data. Identifiers should be kept separate!
- Create "crosswalk" files of identifiers and store them someplace secure.



BOSTON UNIVERSITY

Personal Identifiers


First Name	Last Name	Study ID	Screen ID	Phone #	Identifiable Data
Joseph	Blow	1234	50001	555-131-1111	

Crosswalk file

Subject ID	MRN	SSN
1234	64322	*****

Study Data

ID	Visit	Var 1	Var 2	Var 3
1234	1	5	2	3



HIPAA Identifiers

- Names
- Addresses other than state, and first three digits of the zip code
- All elements of date other than year, and all specific ages over 89 years
- Telephone numbers
- fax numbers
- Email addresses
- Social Security numbers
- Medical Record numbers
- Health plan beneficiary numbers

BOSTON UNIVERSITY

HIPAA Identifiers (cont)

10. Account numbers
11. Certificate/license numbers
12. Vehicle identifiers and serial numbers
13. Device identifiers and serial numbers
14. Web universal resource locators (URLs; web site addresses)
15. Internet protocol (IP) addresses
16. Biometric identifiers, including finger and voice prints
17. Full face photographic images and any comparable images
18. Any other unique identifying number, characteristic, or code



Participant Tracking



Tracking the Participants

You need a system to track participants

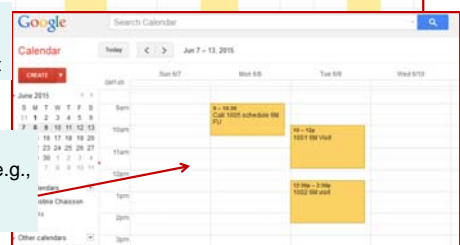
- Tracking for Study Management:
 - Screened, Eligible, Enrolled
 - Monitor and report progress
- Tracking tools for study staff:
 - Schedule/reminders follow up visits
 - Collection of all data points at each visit
- Small study may use Outlook or Excel; large study may need a tracking system



Simple Participant Tracking Tools

ID	Participant Name	phone #	enrollment date	6 M F/U date	6 months F/U target	12 M F/U date	12 M F/U target	18 M F/U date	18 M F/U target
3001	John Smith	617-555-1234	6-Jun-14	6-May-14	6-Jun-14	6-Jul-14	6-Nov-14	6-Feb-15	6-May-15
3002	Jane Doe	617-555-11223	7-Jan-14	7-May-14	7-Jun-14	7-Dec-14	7-Feb-15	7-May-15	7-Aug-15
3003									
3004									

Track in an Excel Spreadsheet



Put into a Calendar (e.g., Google or Outlook)

More Complex Tracking Tool



Tracking the Data Elements

- Identify what data have been collected
 - For each Subject at each Visit:
 - Questionnaires
 - Imaging, labs results
 - Other external data
- Missing data: can you still get it?
- Data cleaning / QA / auditing
- Create "clean" frozen datasets that you append new data to over time



Reports

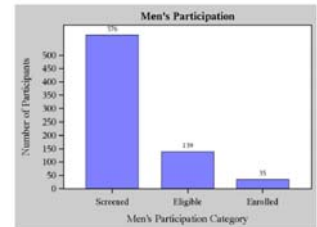
- For the study team to view at regular meetings or online as needed



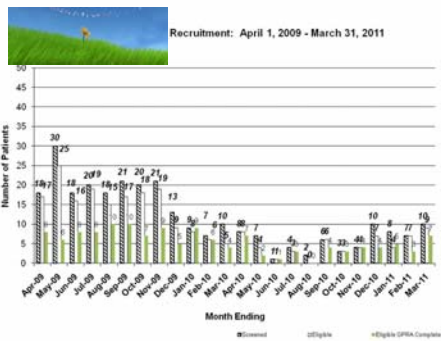
Reports: Visual as well as Tabular

Men's Participation

gender	Status	N	n
MALE:	Screened	576	100
	Eligible	139	24
	Enrolled	35	6



Screened, Eligible and Enrolled



Enrollment

December 2, 2013

Screening and Enrollment Summary

	This Week	Total	HERMAGE	Other Sites
Completed Screening A	215	145	70	
Eligible for Screening B	134 (62.3%)	69 (45.5%)	69 (97.1%)	
Completed Screening B	59 (23.1%)	37 (56.1%)	61 (88.7%)	
Missed Screening B	30 (22.4%)	28 (43.1%)	2 (3.2%)	
Eligible for study	99 (100.0%)	37 (100.0%)	61 (100.0%)	
Enrolled	99 (100.0%)	37 (100.0%)	61 (100.0%)	

Reasons for Ineligibility on Screening A

	n
Ever ART use, Current ART use	70 (46)
Out of catchment area	7
Can't verify HIV status	4
Can't provide fee contacts	4
Can't verify ART naive status	4
Downstage HIV	1
Refusals	1

Baseline Alcohol Use

	Total	Males	Females
Baseline TLFB Complete	99	67	31
Mild or Drinking - Past 7 Days	37 (39%)	22 (33%)	15 (48%)
Mild or Drinking - Past 30 Days	45 (46%)	28 (42%)	17 (55%)
Moderate or Drinking - Past 30 Days	20 (20%)	13 (19%)	7 (23%)
Abstinence - Past 30 Days	32 (34%)	26 (39%)	7 (23%)

Phlebotomy

December 2, 2013

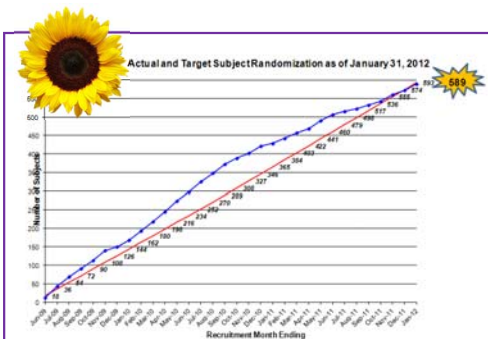
Phlebotomy and Lab Tracking Summary

	n (%)
Blood Drawn - Complete	95 (96%)
Blood Drawn - Incomplete	1 (1%)
Blood Not Drawn 3 attempts	2 (2%)
Blood Not Drawn 2 attempts	1 (1%)
Quality of Venipuncture - Clean	64 (66.7%)
Quality of Venipuncture - Traumatic	22 (23.3%)

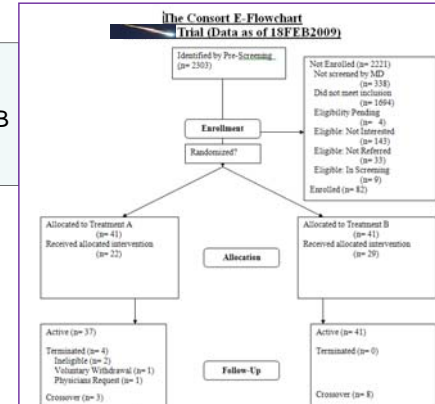
	Mean Time (min)	SD	95% CI	95% ICI
Draw Time (minutes)	2.2 (1.2)	1.0		
Turnout Time (minutes)	12 (2.4)	4.0		
Time to Receipt (minutes)	42 (25.3)	36.3		
Turns Processing, Table 1 (minutes)	64 (28.3)	11.0	60 (26.0%)	67 (28.6%)
Turns Processing, Table 2 (minutes)	64 (28.3)	11.0	61 (26.2%)	67 (28.6%)
Turns Processing, Table 3 (minutes)	59 (26.1)	14.0	42 (34.3%)	69 (37.7%)


Tube	Drawn	Full	Storage Contained	Alcohol Used (0.5mL Used)	Free Cryoprecipitate	CBS Done
Table 1 Purple Top EDTA	64 (97.0%)	67 (94.0%)	63 (89.0%)	6 (7.8)	10	66
Table 2 Blue Top SST	60 (96.0%)	60 (85.0%)	57 (81.0%)	7 (8.2)	10	60
Table 3 Green Top Heparin	62 (96.0%)	61 (87.0%)	51 (84.0%)	10	93 (84.0%)	61 (97.0%)

Actual vs. Targeted Enrollment



Consort Diagram For DSMB Meeting





Participation Summary: Total

Beginning: 5/21/2009 Ending: 4/5/2011

	Number Pending (1)	Number Due (2)	Number Complete (3)	Number Incomplete (4)	Number Inactive (5)	Number Out of Study (6)	Total	Min.	Max.
Baseline	---	---	400	---	0	0	400	---	---
Six Week	0	15	434	1	1	1	450	95.0%	99.3%
Six Month	4	24	395	3	3	0	430	78.2%	98.9%

Participation Summary: ASSIST Score ≥ 4

Beginning: 5/21/2009 Ending: 4/5/2011

	Number Pending (1)	Number Due (2)	Number Complete (3)	Number Incomplete (4)	Number Inactive (5)	Number Out of Study (6)	Total	Min.	Max.
Baseline	---	---	416	---	0	0	416	---	---
Six Week	2	14	393	1	0	0	410	95.1%	99.5%
Six Month	0	24	378	2	3	0	415	78.6%	98.8%

Participation Summary: ASSIST Score < 4




Beginning: 5/21/2009 Ending: 4/5/2011

	Number Pending (1)	Number Due (2)	Number Complete (3)	Number Incomplete (4)	Number Inactive (5)	Number Out of Study (6)	Total	Min.	Max.
Baseline	---	---	44	---	0	0	44	---	---
Six Week	1	1	41	0	1	0	44	95.3%	97.7%
Six Month	4	10	29	0	0	0	43	74.4%	100.0%

(1) Active Subjects who have not yet completed the current interview (not due)
 (2) Active Subjects who are due for the current interview
 (3) Active Subjects who have completed the current interview
 (4) Active Subjects who have missed the current interview within out of their window
 (5) Inactive Subjects
 (6) Out of Study Subjects or Deceased Subjects

Reports

- For the study staff to help manage tasks and know what needs to be done, when

REPORTS

RA Tracking Reports

- Birthday Report
- 6 Week Hair Test Needed
- Contacts Due
- Active Subject Listing

Hair Sample Administrative Reports



- Baseline Hair Sample Status
- Hair Sample Forms Given to DCC
- Invalid Hair Samples
- Missing Hair Samples

IR E Reports

- Scheduled Reporter Sessions
- Missing Reporter Data

Administrative Reports

- All Consented Subjects
- Out of Window --- 6 Weeks
- Out of Window --- 6 Months
- All Missing IRs
- Missed Time Point --- 3 Weeks
- Missed Time Point --- 6 Months

Reports


are logged in to: DCC/Text

- HSN Weekly TEXT Report --- McCard
- HSN Weekly TEXT Report --- St. Mary's
- HSN Weekly TEXT Report --- Mobile Clinics
- HSN Weekly PHONE Report --- McCard
- HSN Weekly PHONE Report --- St. Mary's
- HSN Weekly PHONE Report --- Mobile Clinics
- Intensification Scheme Reminder


3 Month Follow-up Due

- Missing Contact Information
- Missing South African ID Numbers
- HSN Contact Log Lookup
- 3 Month Follow-up Log Lookup
- Missing HSN Enrollment Encounter
- Missing Participant Log Data - TB Treatment --- McCard
- Missing Participant Log Data - TB Treatment --- St. Mary's
- Missing Participant Log Data - TB Treatment --- Clinics
- Missing Participant Log Data - TB Result Obtained --- McCard
- Missing Participant Log Data - TB Result Obtained --- St. Mary's
- Missing Participant Log Data - TB Result Obtained --- Clinics

- 3 Month Follow-up Due
- TB Nurses Follow-up for TB Positive Subjects
- Missing TB Nurse Questionnaire --- McCard
- Missing TB Nurse Questionnaire --- St. Mary's
- Missing TB Nurse Questionnaire --- Clinics
- Missing Participant Log Data - TB Treatment --- McCard
- Missing Participant Log Data - TB Treatment --- St. Mary's
- Missing Participant Log Data - TB Treatment --- Clinics
- Missing Participant Log Data - TB Result Obtained --- McCard
- Missing Participant Log Data - TB Result Obtained --- St. Mary's
- Missing Participant Log Data - TB Result Obtained --- Clinics



Sample Tracking Report: Follow UP




You are logged in to: DCC/Text
 9-Month Follow-Up Due
 Date: 04/04/2013

ID	Enrollment Date	Name	Cell #	Dist #	Family/Friend Dist #	Family/Friend Name	Verification	Target Date	Window Open	Window Close
20737	14/05/2012							14/02/2013	31/01/2013	04/02/2013
10733	14/05/2012							14/02/2013	31/01/2013	04/02/2013
20749	20/05/2012							28/02/2013	14/02/2013	15/04/2013
20756	05/06/2012							05/03/2013	19/02/2013	16/04/2013
20767	15/06/2012							15/03/2013	01/03/2013	26/04/2013
10768	18/06/2012							18/03/2013	04/03/2013	29/04/2013

Look at the Data Early and Often

- You cannot fix a problem if you don't know it exists
- Get data into electronic format that can be manipulated ASAP so it can be more easily reviewed
- Monitor every data point for the first few participants
- Ongoing: audit percentage of forms
- Pay extra attention to key variables



Do simple checks

- Frequency (count) and distribution (range) of each and every variable
- Do crosstabs of variables where appropriate
- What is missing?
- What is out of range?
- What contradicts (e.g., pregnant males)
- Are there systemic problems?



This is why you check...

Scoring Sheet for Kumamoto Scale: [redacted] (as modified by [redacted] August, 2007)

Sex: Male Female
 Height (cm): 165
 Weight (kg): 60.9

Kumamoto Scale
 Yasuma K, Ando Y, Ando C. J. with familial amyloidotic polyneuropathy

I Sensory abnormalities

Lower limbs
 Note the highest level where the copper thermode is felt as cold ≥ 3 of 5 times
 Note the highest level

Upper limbs
 Note the highest level where the copper thermode is felt as cold ≥ 3 of 5 times
 Note the highest level

Trunk and Head
 Note the highest level where the copper thermode is felt as cold ≥ 3 of 5 times
 Note the highest level

II Motor function (muscle)

Ax. Trunk (2: Normal, 3: Normal, 4: Normal, 5: Normal, 6: Normal)
 Quadriceps (2: Normal, 3: Normal, 4: Normal, 5: Normal, 6: Normal)
 Wrist flex (2: Normal, 3: Normal, 4: Normal, 5: Normal, 6: Normal)
 Elbow flex (2: Normal, 3: Normal, 4: Normal, 5: Normal, 6: Normal)

1 SENSORY ABNORMALITIES

Lower limbs
 1. Note the most distal level where the copper thermode is felt as cold ≥ 3 of 5 times
 (1: Toe, 2: Leg, 3: Thigh, 4: not felt at thigh)

2. Note the most distal level where pinsprick is felt ≥ 3 of 5 times
 (1: Toe, 2: Leg, 3: Thigh, 4: not felt at thigh)

3. Note the most distal level where monofilament C is felt ≥ 3 of 5 times
 (1: Toe, 2: Leg, 3: Thigh, 4: not felt at thigh)

Upper limbs
 4. Note the most distal level where the copper thermode is felt as cold ≥ 3 of 5 times
 (1: Finger, 2: Elbow, 3: Shoulder, 4: not felt at shoulder)

5. Note the most distal level where pinsprick is felt ≥ 3 of 5 times
 (1: Finger, 2: Elbow, 3: Shoulder, 4: not felt at shoulder)

6. Note the most distal level where monofilament C is felt ≥ 3 of 5 times
 (1: Finger, 2: Wrist, 3: Shoulder, 4: not felt at shoulder)

Handwritten scores: 4, 4, 4, 2, 3, 4

Medications are never easy

4a-4. For all medications **other than those above**, including non-prescription (i.e. Advil, cold remedies), list the name of your medication in the left-hand column below and click the right-hand columns to indicate which time period(s) you took it.

Medication	1-24 Hours Preceding Cough Attack		25-48 Hours Preceding Cough Attack	
	<input type="radio"/> Yes	<input type="radio"/> No	<input type="radio"/> Yes	<input type="radio"/> No
4a	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
4b	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
4c	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
4d	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
4e	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
4f	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
4g	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
4h	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
4i	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
4j	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Submit

```
data med2=;
  length med2b $10;
  med2b=compress (med2);
  med2c=trim (med2b);
  med2d=lowercase (med2c);
  med2=med2d;

*BASELINE ALLOPURINOL USE-- ghallop;
if gh10txt in ('200mgallup', 'alapernal', 'alapurinol', 'aloeupurin',
'alapurinol', 'allapurino', 'allenpurin', 'allipurano', 'allipurino',
'allipurona', 'allipronol', 'allopurina', 'allopurano', 'allopurina',
'alopurino', 'allopurrin', 'alopureno', 'alourinol', 'allpurinol',
'alluprinol', 'allupurino', 'allupurnoi', 'allupurpin', 'alopurinal',
'alopurinal', 'alpuranal', 'alpurinal', 'aspirin_ib', 'alipurinol',
'blindstudy', 'increasael', 'indomethac', 'juststarte', 'starttedzyl',
'zylorim', 'itookallop', 'allupurina', 'allipurin', 'alpurinol',
'allopurino', 'allpurnal', 'alapurina', 'allopurtno')
then ghallop=1;
```

Perform Systematic Data Audits

- Data forms and source documents are compared with database on X % of forms
- Set an "acceptable" error rate. For example:
 - 0.1% for key variables
 - 0.5% overall
- If audit yields a larger error rate, you must check and correct the database



Audit Example (real data)

6-Month Follow-Up Assessment (Interviewer Administered) - Data Discrepancies

Subject ID	Field Name	CRF	Database	Notes
1115	interdate_6	10/20/08	01/20/2009	Check entire CRF
	Site	1	3	
	Site_other	(text)	-888	
	interstart	12:00	13:30	
	interfinish	12:30	14:00	
	HIV4A_6	Blank	480	
	HIV4A_DK_6	Checked	blank	
	SP3a_1_6	2	1	
	SP4b_6	3	2	
	SP4e_6	15	10	
	SP4f_1_6	0	1	
	SP4f_2_6	0	1	
	SP4f_3_6	0	1	
	SP4g_6	1	-888	
	SP4h_6	1	-888	
	SP4i_1_6	1	-888	
	SP4g_2_6	0	-888	
	SP4g_3_6	0	-888	
	SP4g_4_6	0	-888	
	SP4g_5_6	0	-888	
	SP13_6	5	0	
	SP14_6	1	0	
	SP15_6	2	0	
	SP18_6	1	0	
	STDIG1_6	3	2	

Entered under incorrect ID?

Pay Extra Attention To Key Data

Be sure to pay particular attention to key data points where applicable.

- Query all entries of critical variables (e.g., primary outcome)
- Extra attention to problematic variables (e.g., time-line-follow-back)
- Query all Serious Adverse Events?



Derived Variables

Many analyses require creation of a derived variable from multiple data points

- Be especially careful in creating derived variables
- Include all relevant data elements
- Don't forget to account for missing data
- Be sure to look at frequencies and cross-tabs of derived variables prior to including in models



Creating a Derived Variable

Q1. Does child smoke

Q2. Do household members smoke?

Q3. Do caretakers smoke?

New Var: Smoke_Exp



Q1. Unprotected sex primary partner?

Q2. Unprotected sex with casual partner?

Q3. Share needles?

New Variable: HIV_Exp

Sample SAS Code for Derived Variable

```
if (q1=1) or (q2=1) or (q3=1) then any_exp=1; else
any_exp=2;
```

```
proc freq;
tables any_exp*site;
run;
```

Any_Exp	Yes	No	Total
Site 1	50	40	90
Site 2	20	70	90
Total	70	110	180



Sample SAS Code for Derived Variable

```
/* Corrected code to account for missing */
q1=1) or (q2=1) or (q3=1) then any_expM=1; else
if (q1=0) and (q2=0) and (q3=0) then any_expM=2; else
any_expM=.;
```

```
proc freq;
tables any_expM*site;
run;
```

Any_ExpM	Yes	No	Missing	Total
Site 1	50	40		90
Site 2	20	20	50	90
Total	70	60	50	180



Example 1 missing coded no

The FREQ Procedure

Frequency Percent Row Pct Col Pct	Table of row by col			Total
	row	1	2	
1	50	40	90	
	27.78	22.22	50.00	
	55.56	44.44		
	71.43	36.36		
2	20	70	90	
	11.11	38.89	50.00	
	22.22	77.78		
	28.57	63.64		
Total	70	110	180	
	38.89	61.11	100.00	

Statistics for Table of row by col

Statistic	DF	Value	Prob
Chi-Square		21.0390	<.0001

Example 2 missing coded missing

The FREQ Procedure

Frequency Percent Row Pct Col Pct	Table of row by col			Total
	row	1	2	
1	50	40	90	
	28.46	30.77	59.23	
	55.56	44.44		
	71.43	66.67		
2	20	20	40	
	15.38	15.38	30.77	
	50.00	50.00		
	28.57	33.33		
Total	70	60	130	
	53.85	46.15	100.00	

Statistics for Table of row by col

Statistic	DF	Value	Prob
Chi-Square	1	0.3429	0.5576

What's up with the missing values?

- Go back and look at forms:
 - Is there an explanation?
 - Is the missing data differential?
- What are the implications?
 - Example: There are 2 sites and all the forms with missing values came from a single site
- Did you find this problem early enough to correct it?
- This is why you check "early and often"



Data Security - General

- Keep paper records should be kept in locked cabinets and/or offices
- Store identifiers like names and addresses separate from clinical data
- Keep particularly sensitive data apart from other identifiers (e.g., SSN) – in a separate file, by ID
- Do not collect sensitive data unless you really need it



Data Security - Hardware

- Password protect all computers
- Set to automatically timeout if inactive
- Encrypt laptops, flash-drives and other storage devices when possible
- Do not put identifiable data on portable media (e.g., CDs, flash-drives) unless password protected, preferably encrypted



Take Home Message

- Your team should include someone who understands data issues
- Budget for data management
- Planning ahead results in fewer revisions
- Check your data early and often
- If you do things correctly from the beginning:
 - It is less work
 - It is less expensive
 - You are more likely to discover the truth at the end



Questions?

