## DATA MANAGEMENT

What you should know and why you should care

CREST Seminar
November 10, 2015

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

**BU**

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 be 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

2

---

## Real World Examples

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

3

---

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**

---

### Oops: Excel Error Calls Into Question…
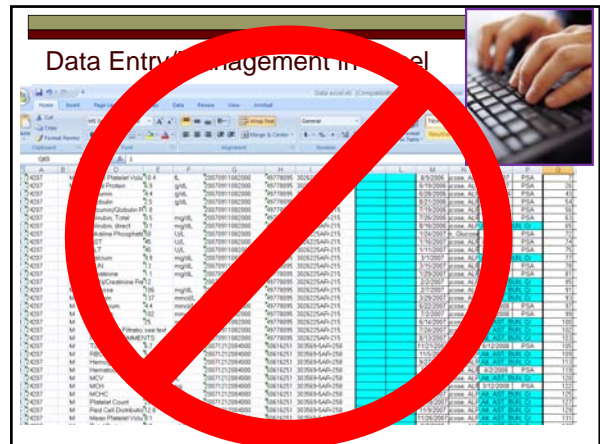
IEEE **SPECTRUM**        Posted 22 Apr 2013

- A book by Harvard Researchers entitled *'This time it's different"* contained "…serious errors that inaccurately represent the relationship between public debt and GDP growth among 20 advanced economies in the post-war period."
- The Authors admitted they forgot to include five rows in an Excel file resulting in exclusion of data from 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, 2011

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



## Data Management 101:

- 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 manage and share your data; meet funder requirements. General elements include:

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

## 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

---

## Timeline



---

## 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)

---

## Other Data Elements

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

---

## 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

---

## Sample data/visit grid

## 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 | 6-12 | 12-18 | 18-24 | 24-30 | 30-36 | 36-42 | 42-48 | 48-54 | 54-60 |
|---|---|---|---|---|---|---|---|---|---|---|
| Hiring & training | | | | | | | | | | |
| Finalize instruments & IRB | | | | | | | | | | |
| Enrollment | X | X | X | X | | | | | | |
| Intervention | X | X | X | X | X | X | X | | | |
| Follow-up | | | X | X | X | X | X | X | | |
| Data QA/clean | | | X | X | X | X | X | X | X | |
| Primary & secondary data analyses | | | | | | | | | X | X |
| Presentations & Publication | | | | | | | | | X | X |
| Study meetings | X | X | X | X | X | X | X | X | X | X |

Design/construct data collection systems

## Sample Task-based Gantt

| | Year 1 | | | | Year 2 | | | |
|---|---|---|---|---|---|---|---|---|
| Tasks | Months 1-3 | Months 4-6 | Months 7-9 | Months 10-12 | Months 13-15 | Months 15-18 | Months 29-21 | 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 | | | | | | | | |
| DSMB data freeze, reports, meeting | | | | | | | | |
| Follow up visits | | | | | | | | |

## Multi-task Indicating Responsible Parties



## 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

## Sample Data Dictionary: SAS formatted dataset

## 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

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

Data Manager → Case Report Form ← Study Coordinator

Co-Investigators → Case Report Form ← Study Sponsor?

Biostatistician → Case Report Form ← Clinicians
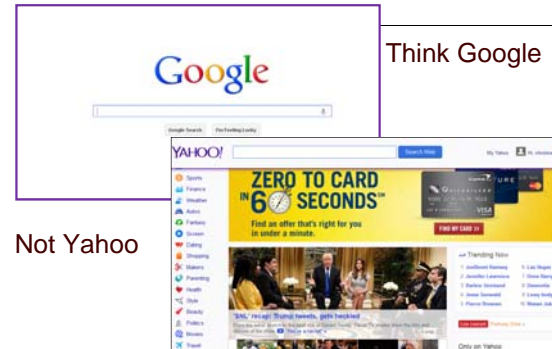
## Why is this topic important?

- 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…

## Successful Form: Consider ALL Functions

- ❑ Data Collection - who is completing form?
  - ❑ Study Staff (Coordinator, Clinician)
  - ❑ Participant
  - ❑ Clinician
- ❑ Data entry - who is entering data?
  - ❑ Study staff
  - ❑ Outsourced
- ❑ Data management/cleaning
- ❑ Auditing
- ❑ Data analysis

---

When designing forms…



Think Google

Not Yahoo

---

## What makes a good 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

---

## Pilot Your Forms *Prior* to Data Collection

- ❑ Test in target population (e.g., age, literacy)
  - ■ Are items left blank?
    - ➤ Reword/drop question
  - ■ Are "skip" patterns followed correctly?
    - ➤ Train clinic personnel/revise or simplify forms
  - ■ Are open-ended questions generating common responses?
    - ➤ Categorize/code
- ❑ Make corrections prior to start of study
- ❑ Do not start data collection until forms are final

---

## Avoid Open-ended & Include Response Measure

| | |
|---|---|
| What is your date of birth? _____ <br><br> How much do you weigh? ____ <br><br> How tall are you? _____ <br><br> Record subject's temperature _____ | Date of Birth? __/__/____ <br>          MM DD YYYY <br> How much do you weigh? __.__ <br>             (pounds) <br> How tall are you? ___ /__ <br>          (feet)/(inches) <br><br> Record subject's <br> oral temperature   __ __. __ (f) |

---

## Include Clear Instructions

A. What is your race/ethnicity? **(Check one)**
- $_1$O Caucasian
- $_2$O African American/Black
- $_3$O Asian, Pacific Islander
- $_4$O Native American
- $_5$O Other _____

B. What is your race/ethnicity? **(Check all that apply)**
- $_1$❑ Caucasian
- $_1$❑ African American/Black
- $_1$❑ Asian, Pacific Islander
- $_1$❑ Native American
- $_1$❑ Other _____

## Beware of "check all that apply"

| | COMORBIDITIES | |
|---|---|---|
| b. | Heart disease | ₁☐ |
| c. | Diabetes | ₁☐ |
| d. | Hypertension | ₁☐ |
| e. | Pulmonary disease | ₁☐ |

| a. | COMORBIDITIES | |
|---|---|---|
| b. | Heart disease | ₁○ ₂○ |
| c. | Diabetes | ₁○ ₂○ |
| d. | Hypertension | ₁○ ₂○ |
| e. | Pulmonary disease | ₁○ ₂○ |

37

---

## Account For Missing Data

| CBC | Unit | Value | |
|---|---|---|---|
| 1. Hemoglobin | g/dl | __ __.__ | ☐ Not Done |
| 2. Hematocrit | % | __ __.__ | ☐ Not Done |
| 3. RBC | M/mm³ | __ __.__ | ☐ Not Done |

---

## Specify the Units



Alcoholic Beverages
Serving Sizes

A 12-ounce bottle or can of beer = 1 serving

A 1-1.5 ounce shot of liquor straight or in a mixed drink = 1 serving

A 5-ounce glass of wine = 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          ☐₁ USA

☐₂ Guatemala    ☐₂ Guatemala

☐₃ Mexico       ☐₃ Mexico

☐₄ Dominican Republic    ☐₄ Dominican Republic

☐₅ Other [____]    ☐₆ El Salvador

☐₅ Other

---

## ID Assignment

☐ Must be UNIQUE for each subject

☐ Should appear on every form (preferably page)
- Links paper form with specific record in database
- Multiple forms, "merge key" in database

☐ 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

Do NOT do this



Overdose Prevention & Naloxone Form

Codes and Abbreviations:
FtM     Female to Male transgender
MtF     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



## Consider who is Completing a Form



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" sofware packages available
- No single best choice
- Cost can vary widely
- Database structures vary
- Do your homework to make sure what you get will work for your project

---



Find out what software is available through your institution

---



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)

53

## 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!)

## Crosswalk for Personal Identifiers

- 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.

## Personal Identifiers

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

Identifiable Data

| Subject ID | MRN | SSN |
|---|---|---|
| 1234 | 64322 | ****** |

Crosswalk

Study Data

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

I ♥ my Privacy

## HIPAA Identifiers

1. Names
2. Addresses other than state, and first three digits of the zip code
3. All elements of date other than year, and all specific ages over 89 years
4. Telephone numbers
5. fax numbers
6. Email addresses
7. Social Security numbers
8. Medical Record numbers
9. Health plan beneficiary numbers

## 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



I think we really need to review our tracking policy, Mr.Corello!

SIPRESS

## 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



Track in an Excel Spreadsheet

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

## More Complex Tracking Tool



## Reports

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



Consort Diagram



## Enrollment          Phlebotomy



## Reports: Visual as well as Tabular

## Actual vs. Targeted Enrollment



## Screened, Eligible and Enrolled





## Reports

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

## Sample Tracking Report: Follow UP



## 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 entry: 1st entry/ 2nd entry/ reconciled
- Data cleaning/ QA / auditing
- "Clean" frozen datasets

## Look at the Data Early and Often

- You cannot fix a problem if you don't know it exists
- Get data into electronic format 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…



## Medications are never easy

```
data meds1;
    length med2b $10;
    med2b=compress(med2);
    med2c=trim(med2b);
    med2d=lowcase(med2c);
    med2=med2d;

*BASELINE ALLOPURINAL USE-- ghallop;
if gh10txt in ('200mgallup' 'alapernal' 'alapurinol' 'aleenpurin'
'alepurinol' 'allapurino' 'allenpurin' 'allipurano' 'allipurino'
'allipurona' 'allipronol' 'allopruina' 'allopurano' 'allopurina'
'allopurino' 'allopurrin' 'aliopureno' 'allourinol' 'allpurinol'
'alluprinol' 'allupurino' 'allupurnoi' 'allurpurin' 'alopurinal'
'alopurinol' 'alpuranal' 'alpurinal' 'aspirin,ib' 'alipurinol'
'blindstudy' 'increaseal' 'indomethac' 'juststarte' 'startedzyl'
'zyloprim'   'itookallop' 'allupurina' 'allinpurin' 'alpurinol'
'allopurino' 'allpurnal'  'alapournia' '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 | 03/30/2009 | 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 | Q1. Unprotected sex primary partner? |
| Q2. Do household members smoke? | Q2. Unprotected sex with casual partner? |
| Q3. Do caretakers smoke? | Q3. Share needles? |
| New Var: Smoke_Exp | New Variable: HIV_Exp |

## Slide 1

### 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 |

## Slide 2

### 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=.; else
        any_expM=2;

  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 |

## Slide 3



Example 1 missing coded no / Example 2 missing coded missing — The FREQ Procedure tables and Chi-Square statistics.

## Slide 4

### 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"

## Slide 5

### A Caution About New Technology



Too soon, not ready for prime time

Too late, becoming obsolete

## Slide 6

### Smartphones

## Benefits of smartphones

- Electronic data capture
- Secure if you use to connect to website with encryption
- SMS (text messaging) in addition to linking to web-form
- Easy to carry
- Everyone wants one
- "Sexy" so funders like the idea

## Smartphone challenges

- Can be Expensive (hardware, data plans)
- Cannot encryption text messages
- One question per screen
- Small screens make view some question types difficult (e.g., grids)
- Navigating around questionnaire (going back) is challenging
- Battery life is short
- Attractive to thieves and easily stolen

## 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?