

# Preventive Medicine Grant Writing Seminar Series

## Session 7: Data Collection, Data Management, Storage & Quality Assurance

*Phyllis A. Richey, Ph.D.*

*Professor, Departments of Preventive Medicine, Physical Therapy and Pediatrics  
Clinical & Community Research Informatics Section (CCRIS)*



2

## Overview

- **Browner et al chapters 18-19**
- **Informatics vs Research Informatics**
- **Data Collection, Management Defined Entry Systems**
  - Database design for study implementation, participant recruitment and retention, outcome measure data collection, intervention delivery
  - Data standardization
- **Technology Resources**
  - Where used
  - How processed
- **Programming, Data Entry, Cleaning, Analysis & Security**
- **Special insights on research data management design**

Preventive Medicine Grant Writing Seminar Series: Session 7



3

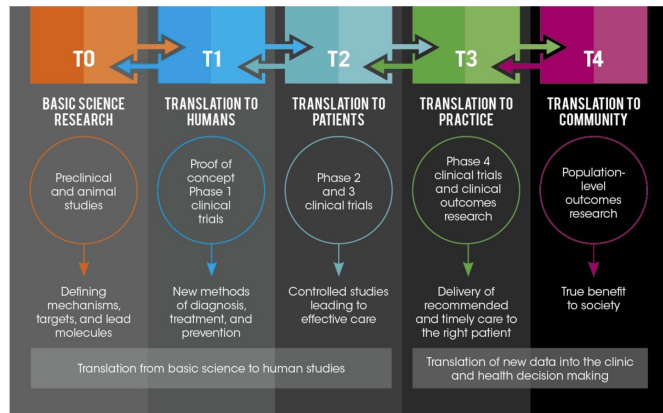
# Informatics vs. Research Informatics

## ● INFORMATICS

- The science of processing data for storage and retrieval.

## ● RESEARCH INFORMATICS

- Sub-discipline within biomedical informatics which focuses on developing new theories, tools, and solutions used to translate data across the basic research, clinical trials, medical center and community practice continuum.



\*National Center for Advancing Translational Sciences (NCATS), 2023

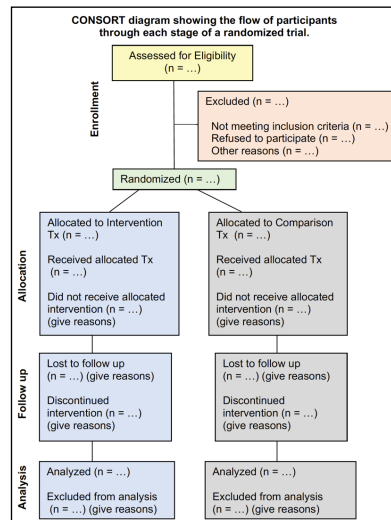
## Why is a research informatics solution critical to successfully conducting a research study?

- **It defines the procedures for information (data) processing and storage during the study**
  - Determines the methods of data collection and storage
  - Establishes user interface (GUI) through which information is entered, navigated, and presented in a meaningful way.
  - Provides manipulation of data and creation of new information through rules, calculations, index lookups

## Development of a research informatics solution:

- Defines “HOW” the study is to be conducted

- Participant Recruitment and Enrollment - Where the CONSORT begins
- Eligibility and Randomization
- Intervention Delivery
- Data Collection – outcome and process measures
- Storage & Mining of Study Information
- Data Sharing (central repository)
- Ongoing Quality Assurance for Data Integrity
- Data Security Maintenance – throughout the life of the project



## Research Informatics Solutions

- Common Types of Data Management Resources
  - Single spreadsheet
  - Statistical
  - Database Management System (DMBS)
- The Foundational Unit – The Data Table

# So ... What's a data table?

- Organized representation of data into columns and rows
- Data organization is strictly restricted
  - One variable (field) per column
  - One record per row

# Example of what is NOT a Data Table

The screenshot shows an Excel spreadsheet titled 'YZSORT IMPACT Data - Compatibility Mode - Excel'. The data is organized into several sections: 'District Report Data', 'School Report Data', and 'Average' data. The columns include 'SITE', 'SCHOOL', 'PES', 'Barcode ID', 'POS#', and various activity categories like 'Lying', 'Sitting', 'Standing', 'Walking', 'Very Active', 'MVPA', 'Management', 'General Knowl', 'Physical Fitne', 'Fitness Activi', and 'Fitness Activity'. The data is presented in a way that does not strictly follow the 'one variable per column' rule, as some columns contain multiple variables or are grouped together.

## Example of a TRUE Data Table

Record ID	RZID	LNAME	FNAME	Street	City	State	Zip
1	001	Jones	Joe	1900 Poplar Ave	Memphis	TN	38105
2	002	Smith	Mary	9215 Oak Rd	Little Rock	AR	72002
3	003	Green	John	Stage Rd	Bartlett	TN	38134
4	004	Allen	Tom	1 Avenue	Jackson	MS	38108
5	005	Mouse	Micky	1 Main Street	Orlando	FL	32801
6	006	Duck	Donald	2 Main Street	Orlando	FL	32802
7	007	Kent	Clark	100 Hickory Lane	Smallville	KS	66605

## Data Management Resource Types

- **Single spreadsheet**
  - Most basic - Limited to two-dimensional data table format
- **Relational Database**
  - Allows for more complex data format
- **Each data management resource type's unique features = its "Data Management Model"**

ParticipantID	FName	DOB	Sex	Hyperbili_ind	ExDate	ExWght	ExHght	IQ
2101	Robert	1/6/2010	M	1	1/29/2015	23.9	118	104
2322	Helen	1/6/2010	F	0	1/29/2015	18.3	109	94
2376	Amy	1/13/2010	F	1	3/22/2015	18.5	117	85
2390	Alejandro	1/14/2010	M	0				
2497	Isiah	1/18/2010	M	0	2/18/2015	20.5	121	74
2569	Joshua	1/23/2010	M	1	2/13/2015	24.8	113	115
2819	Ryan	1/26/2010	M	0				
3019	Morgan	1/29/2010	F	0	2/9/2015	19.1	105	105
3031	Cody	2/15/2010	M	0	4/16/2015	15.2	107	132
3290	Amy	2/16/2010	F	1	4/12/2015	18.0	102	125
3374	Zachary	2/21/2010	M	1				
3625	David	2/22/2010	M	1	2/10/2015	19.2	114	134
3901	Jackson	2/28/2010	M	0				

■ **FIGURE 19.1** Simplified data table in "datasheet view" for a cohort study of the association between neonatal hyperbilirubinemia and IQ score at age 5. The binary predictor is "Hyperbili\_ind," defined as whether the total bilirubin rose to 25 mg/dL or more in the first 10 days after birth, and the continuous outcome is "IQ," the participant's IQ score at age 5. Participants 2390, 2819, 3374, and 3901 were not examined at age 5.

■ **FIGURE 19.2** A two-table infant jaundice study database with a table of study participants (in which each row corresponds to a single study participant) and a table of examinations (in which each row corresponds to a particular examination). For example, Participant 2322 is identified as Helen, date of birth 1/6/2010, in the first table, and has three exams in the anonymous second table. Note that ExWght and ExHght are entered in the exam table, not the participant table.

(Browner et al, Designing Clinical Research 5th Edition, 2023)

# Advantages of Database Management System (DBMS) Modeling

- **What is Data Modeling?**

- Techniques, not rules
- Real world problem translated to a digital structure
- Good modeling is essential

## Types of DBMS Models

- **Depends on how the data is organized**

- **Flat Design Model**

- Presents data in a single table or list (address list)
- Flat design models have a “One to One Relationship”
- Fields represent all parameters (prone to duplicate data)
- Functionality limited to storing information, manipulate fields, printing and displaying data

	A	B	C	D	E	F	G	H
1	Record ID	RZID	LNAME	FNAME	Street	City	State	Zip
2	1	001	Jones	Joe	1900 Poplar Ave	Memphis	TN	38105
3	2	002	Smith	Mary	9215 Oak Rd	Little Rock	AR	72002
4	3	003	Green	John	Stage Rd	Bartlett	TN	38134
5	4	004	Allen	Tom	1 Avenue	Jackson	MS	38108
6	5	005	Mouse	Micky	1 Main Street	Orlando	FL	32801
7	6	006	Duck	Donald	2 Main Street	Orlando	FL	32802
8	7	007	Kent	Clark	100 Hickory Lane	Smallville	KS	66605

# Types of DBMS Models

## ● Relational Design Model

- Developed by E.F. Codd in 1970 (Ref: “A Relational Model of Data for large Shared Data Banks”)
- Excellent for data entry
- Utilizes multiple (flat) tables - single data points per field (normalized)
- Eliminates redundancy through “normalization” (smoothes out data clumps)
- Functionality to perform queries, join tables, establish integrity constraints, create reports
- Utilizes unique fields defined as indexes (keys) establishing connections (relationships) between tables
- Useful in sharing data across networks, internet, electronic devices, other (software) programs
- Relational Models can have One to One, One to Many, Many to One or Many to Many Relationships

# One versus Many

## Defining Different Types of Relationships

### ● One to one (typical flat table)

- Person and Social Security Number
- Patient and blood type
- Rare to see these in separate tables

	A	B	C
1	<b>First Name</b>	<b>Last Name</b>	<b>Blood Type</b>
2	Mary	Green	A+
3	John	Smith	A+
4	Glenda	Robbinson	O+
5	Susan	Jones	B-
6	Bobby	Taylor	A-
7	Brian	Williams	B+
8	Fredrick	Miller	O-
9	Alice	Washington	B+

# One versus Many

## Defining Different Types of Relationships

- **One to many / many to one**

- Team to player
- Blood samples to patient

Patients

Patient ID	First Name	Last Name	Blood Type
001	Mary	Green	A+
002	John	Smith	A+
003	Glenda	Robbinson	O+
004	Susan	Jones	B-
005	Bobby	Taylor	A-
006	Brian	Williams	B+
007	Fredrick	Miller	O-
008	Alice	Washington	B+
009	Kelly	Franklin	A+
010	Bill	Davis	O+

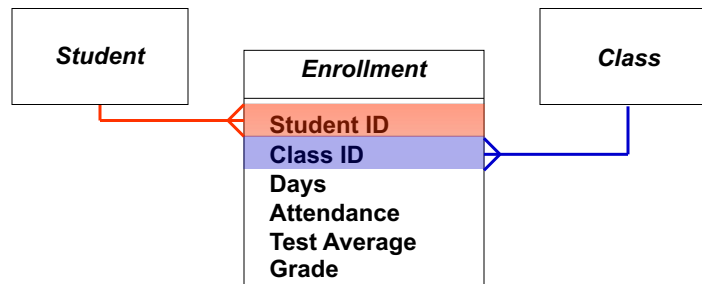
Blood Samples

Sample ID	Patient ID	Date	Time
1	003	10/25/21	6:00
2	004	1/5/22	12:16
3	001	12/5/19	13:56
4	002	1/2/20	20:03
5	005	11/16/20	6:30
6	008	8/8/22	6:42
7	001	4/3/22	16:01
8	006	11/1/21	9:14
9	003	2/19/21	23:00
10	004	7/7/19	18:37
11	006	11/4/20	3:15
12	007	3/5/21	14:59
13	008	3/5/21	10:45
14	009	2/5/22	11:07
15	010	9/5/21	12:45
16	006	12/15/19	12:05
17	004	4/5/22	14:50
18	005	6/18/20	15:19
19	006	1/25/22	16:45
20	007	7/5/19	17:49

(Browner et al, Designing Clinical Research 5th Edition, 2023)

## Modeling: Many to Many Relationships

- How do we model these relationships?
- Two one to many relationships
- Can't just use one key (connector) between these 2 files
- We need a third table "in the middle"

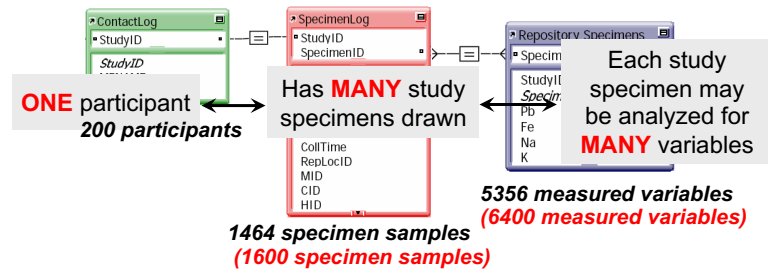




# Many to Many Relationships

## Role in Efficient Data Storage

- The most complex relationship type
  - Example:
    - Patient blood samples measured for concurrent medications



- Generally, need to do extra work to resolve these

# Common Methods of Data Collection & Storage

Software Programs Used in Database Design, Data Management & Analysis					
Spreadsheet	Integrated Desktop Relational Database	Enterprise Database	Integrated Web-Based Data Mgmt	Machine Readable Forms & Online Survey	Statistical Analysis
MS Excel	FileMaker Desktop, Server, App* (Cross platform)	Oracle	REDCap*	Qualtrics	R*
Google Drive Spreadsheet (Sheets*)	MS Access Office 365 - Desktop (PC only)	SQL	MediData RAVE	QuestionPro	SAS
Apache OpenOffice*		MySQL	EpilInfo*	TeleForm	SPSS
		PostgreSQL*	Datalabs EDC	Zoomerang	Stata
			QuesGen	SurveyMonkey	JMP
			OnCore		

\* Free

(Partially reproduced from Browner et al, Designing Clinical Research 5th Edition, 2023, Table 19.1)

# DBMS Programming Components

- **Data Tables**
- **Data Modeling Relationships**
  - Primary and Foreign Keys
  - Normalization
- **Data Dictionaries (Metadata)**
  - Field name (coding conventions)
  - Field data type (text, number, date, time, container)
  - Field definition
- **Field validation constraints**
  - Auto-entered, calculated, looked-up values
  - Restricted data type, not empty/unique, within range, value list
- **Data entry control style**
  - Edit box, value pick list, radio button, check box, calendar

# DBMS Data Entry Methods

- **Case report form (paper form with manual transcription)**
- **OCR (Optical Character Recognition) - TeleForm**
- **Direct electronic data capture**
  - Internal DBMS layouts – desktop, laptop, mobile device (FileMaker Go App)
- **Indirect data capture**
  - External data import protocols – Open Database Connectivity (ODBC), Java Database Connectivity (JDBC)
    - Survey tools – Qualtrics, QuestionPro, REDCap, SurveyMonkey
    - Technology devices – Accelerometers, BP/Holter monitors, body weight scales, lab equipment
    - Apps – Loselt, MyFitnessPal
    - Online Validated Instruments – NCI's Diet History Questionnaire (DHQ)

## DBMS Data Cleaning, Analysis and Confidentiality

- **Queries and Reports**
  - Filter data using queries
  - Check for errors
  - Summarize data in reports
- **Statistical Analysis**
  - Extract “cleaned” datasets
- **De-identification & Resource Sharing**

## DBMS Security

- **Data access and processing requires network authentication**
- **Server Hosted Database Solutions - use registered TCP/IP port**
  - Users authenticated at server-level, not client-level
  - Encrypted network traffic - Server uses TripleDES encryption with HMACSHA-1 for integrity checking (168 bit key)
  - Additional 128 bit SSL between FileMaker server and client
- **VPN is required if accessing from outside of network**
- **e-Data entry verified via record-level audit log**
- **Unique usernames with passwords required**
- **Logon/logoff whenever switching users**



# Questions?

