



Collecting rich paradata to monitor data collection quality

Gina-Qian Cheung & Beth-Ellen Pennell
University of Michigan

3MC International Conference
July 27th, 2016
Chicago



Overview

- Paradata collection
- Paradata usage examples
- Some observations



Paradata includes...

- Interviewer (experience, training grades, historical performance)
- Sample segments (PSU, Stratum, observations)
- Address (probability of selection, observations, # contacts, status)
- Screener contacts (call #, interviewer, time, date, informant behavior, outcome)
- Household (composition, informant behavior, sample respondent characteristics)
- Main interview contacts (call #, interviewer, time, date, informant behavior, outcome)



Paradata also includes...

- Audit trails
 - Screener and survey interview (keystrokes, timings, functions, consistency checks, suspensions)
 - Sample management system (log and timing of actions)
- Digital photos
- Fingerprints
- GPS (Global Positioning System)
- Digital recordings
- Collection of various anthropometric data using digital devices



User Agent String for Web Users

Most Web browsers use a User-Agent string value as follows:

Mozilla/[version] ([system and browser information]) [platform] ([platform details]) [extensions].

For example, Safari on the iPad has used the following:

```
Mozilla/5.0 (iPad; U; CPU OS 3_2_1 like Mac OS X; en-us)  
AppleWebKit/531.21.10 (KHTML, like Gecko) Mobile/7B405
```

The components of this string are as follows:

- *Mozilla/5.0*: Previously used to indicate compatibility with the Mozilla rendering engine.
- *(iPad; U; CPU OS 3_2_1 like Mac OS X; en-us)*: Details of the system in which the browser is running.
- *AppleWebKit/531.21.10*: The platform the browser uses.
- *(KHTML, like Gecko)*: Browser platform details.
- *Mobile/7B405*: This is used by the browser to indicate specific enhancements that are available directly in the browser or through third parties



Two Examples

- Ghana Socioeconomic Panel Study
- The China Health and Retirement Longitudinal Study



Ghana Socioeconomic Panel Study

- Revisit panel households at 3-4 year intervals for 20 years.
 - Sponsored by Economic Growth Center (EGC) at Yale University
 - Carried out by the Institute for Statistical, Social and Economic Research (ISSER) at the University of Ghana.
- First wave (baseline) was completed on ***paper*** between October 2009 and February 2010.
- Second wave was conducted on ***Computer-Assisted Personal Interviewing (CAPI)*** between March-December 2014.
 - Collaborated with Survey Research Center (SRC) at University of Michigan.
- 334 enumeration areas country-wide. Sample size of 5009 households, with approximately 18,000 individuals. Also sample size of 500 split-off households were tracked and interviewed between January-June 2015.
- Interviews are **NOT digital recorded** for quality monitoring purpose

Instrument design

- Interviewers have a high-level of autonomy with respect to interview navigation.
- Interviewers are able to:
 - switch respondents easily.
 - jump to any section of questionnaire quickly.
- Development of a questionnaire “Dashboard” to show the status of all the questionnaire sections and all the respondents within the household.

Using Paradata for Questionnaire Design

- How does instrument design affect instrument navigation?
 - **Instrument parallel blocks**: four instruments (household, personal, agriculture, enterprise) with multiple sections/blocks within each instrument.
- How does instrument navigation affect interview length?
 - **Order of interview initiation**
 - **Movements between blocks**

By using keystroke data (Paradata)

Timestamps

Hours:Minutes:Seconds:Thousands of a second

Case ID in Blaise database

"1/17/2012 9:00:06:304 AM", "Enter Form:1", "Key:3975053020 " ← **Sample ID**

"1/17/2012 9:00:06:304 AM", "Metafile name:C:\blproj\HRS2012\work\HRS12.bmi" ← **Start IW**

"1/17/2012 9:00:06:304 AM", "Metafile timestamp:Friday, January 06, 2012 1:08:04 PM" ← **Audit trail file information**

"1/17/2012 9:00:06:304 AM", "WinUserName:14554015" ← **Interviewer ID**

"1/17/2012 9:00:06:304 AM", "DictionaryVersionInfo:0.0.0.0"

...

"1/17/2012 9:00:12:702 AM", "Enter Field:SecA.StartInterview.A007TRALive_A", "Status:Normal", "Value:"

"1/17/2012 9:00:13:965 AM", "(KEY:)1[ENTR]" ← **Time of first keystroke**

"1/17/2012 9:00:14:276 AM", "Action:Store Field Data", "Field:SecA.StartInterview.A007TRALive_A" ← **Question**

"1/17/2012 9:00:14:328 AM", "Leave Field:SecA.StartInterview.A007TRALive_A", "Cause:Next Field", "Status:Normal", "Value:1"

...

"1/17/2012 9:02:51:681 AM", "Enter Field:SecJ.WORKSTATUS.J005MCurrEmpStatus[1]", "Status:Normal", "Value:"

"1/17/2012 9:02:55:971 AM", "(KEY:)15[BACK][BACK]5[ENTR]" ← **Question with change answer**

"1/17/2012 9:03:03:209 AM", "Action:Store Field Data", "Field:SecJ.WORKSTATUS.J005MCurrEmpStatus[1]"

"1/17/2012 9:03:03:256 AM", "Leave Field:SecJ.WORKSTATUS.J005MCurrEmpStatus[1]", "Cause:Next Field", "Status:Normal", "Value:5"

...

"1/17/2012 9:13:24:923 AM", "Enter Field:IWComplete", "Status:Normal", "Value:"

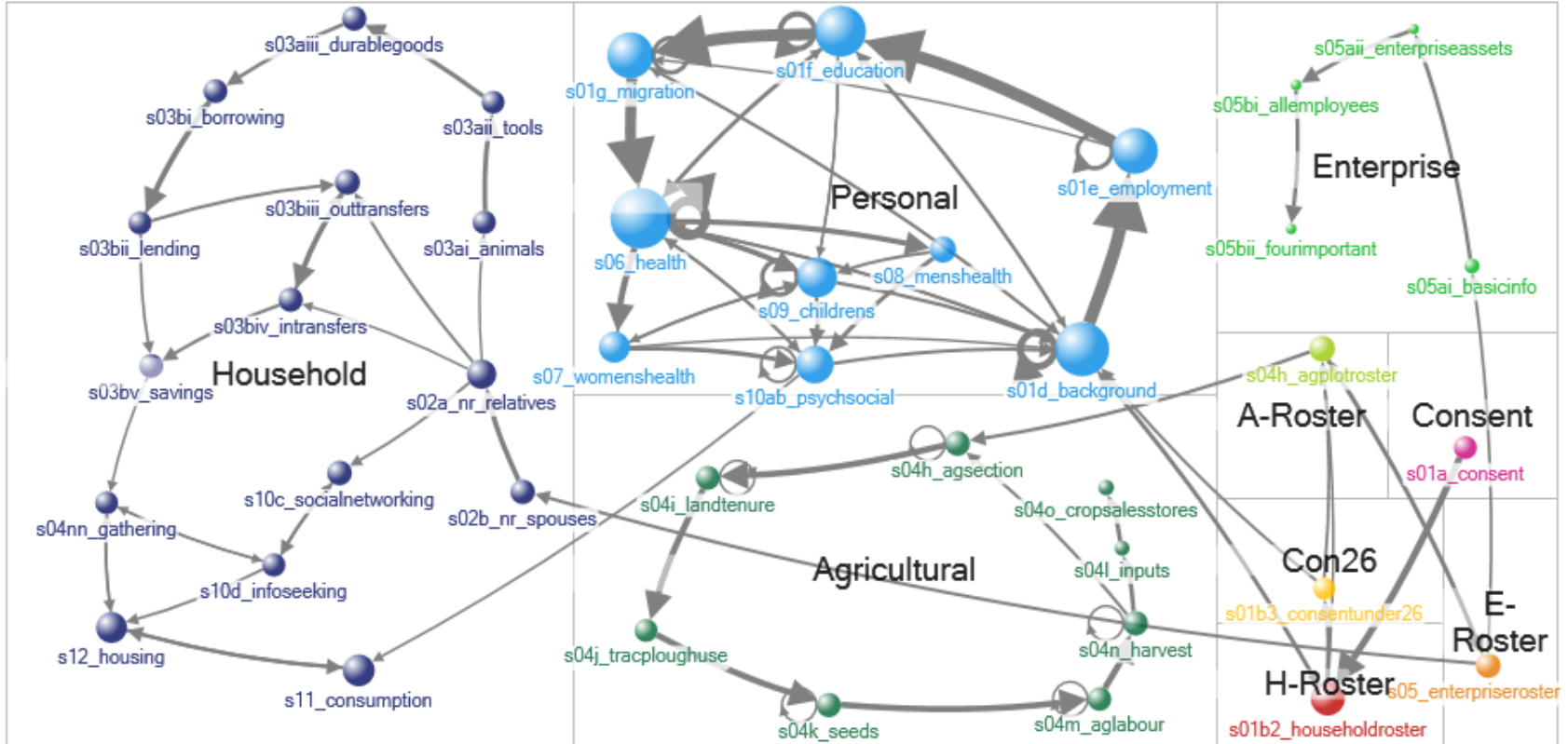
"1/17/2012 9:13:28:480 AM", "(KEY:)1[ENTR]"

"1/17/2012 9:13:29:650 AM", "Action:Store Field Data", "Field:IWComplete"

"1/17/2012 9:13:29:728 AM", "Leave Field:IWComplete", "Cause:Next Field", "Status:Normal", "Value:1" ← **Complete IW**

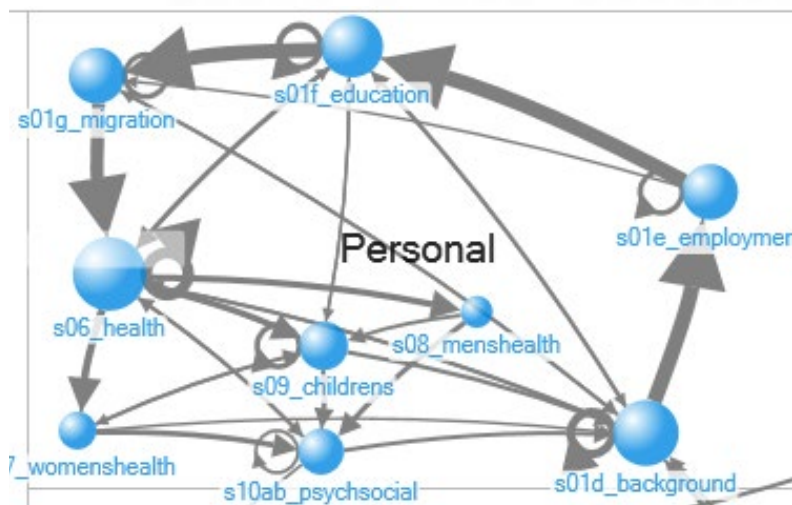
"1/17/2012 9:13:30:056 AM", "Leave Field:IWComplete", "Cause:Exit", "Status:Normal", "Value:1"

"1/17/2012 9:13:30:056 AM", "Leave Form:1", "Key:3975053020 "



Most Common Block Moves All Types

- Edge Weight (number of times a move occurred) ≥ 500
- Movement within sections dominates
- Exceptions are rosters and Personal to Household



Blaise 4.8 Data Entry - c:\blproj\ghana_p\work\householdsurvey

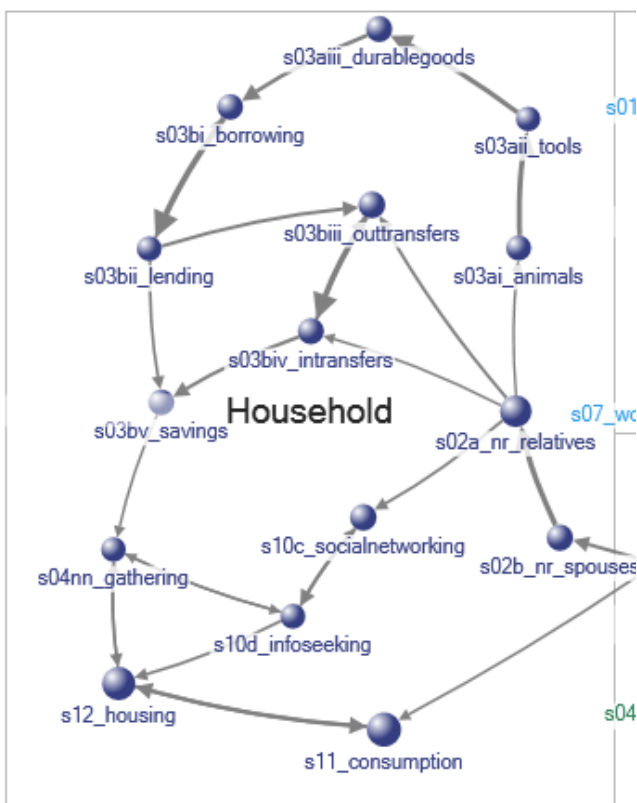
Forms Answer Help

HOUSEHOLD SURVEY Person Status Enterprise Status Agriculture

Name	Background	Employment	Education	Migration	Health	Womens Health	Mens Health	Children	Pysch/Social
ADAM K (AK)	Started	Done	Done	Done	Started	---n/a---	Not Started	---n/a---	Not Started
AMINA A (MINA)	Not Started	Not Started	Not Started	Not Started	Not Started	Not Started	---n/a---	---n/a---	Not Started
ABDUL A	Not Started	Not Started	Not Started	Not Started	Not Started	---n/a---	Not Started	---n/a---	Not Started
TANLIDOW A	Not Started	Not Started	Not Started	Not Started	Not Started	---n/a---	---n/a---	Not Started	---n/a---
YUSSIF A	Done	---n/a---	Not Started	---n/a---	Not Started	---n/a---	---n/a---	Not Started	---n/a---
LAILATU A	Done	---n/a---	---n/a---	---n/a---	Not Started	---n/a---	---n/a---	Not Started	---n/a---

Most Common Block Moves All Types

- Tendency to move laterally or within the same questionnaire content
- Optional sections introduce multiple, common paths



HOUSEHOLD SURVEY

Person Status

Enterprise Status

Agriculture

Survey Status

S0Consent Forms

S01A: Consent

Done

S01B3: Consent for Under 26

Done

Rosters

S01B2: Household Roster

Done

S04: Plot Roster

Done

S05: Non-Farm Enterprise Roster

Done

Person Sections

Started

Plot Sections

Done

Non-Farm Enterprise Sections

Not Started

Household Level Sections

S02B: Non-Resident Spouses

Not Started

S02A: Non-Resident Relatives

Not Started

S10C: Social Networking

Started

S10D: Information Seeking

Done

S11: Household Consumption

Not Started

S12: Housing Characteristics

Not Started

S04NN: Gathering

Done

S03Ai: Animals

Done

S03Aii: Tools

Started

S03Aiii: Durable Goods

Started

S03Bi: Borrowing

Not Started

S03Bii: Lending

Done

S03Biii: Out-Transfers

---n/a---

S03Biv: In-Transfers

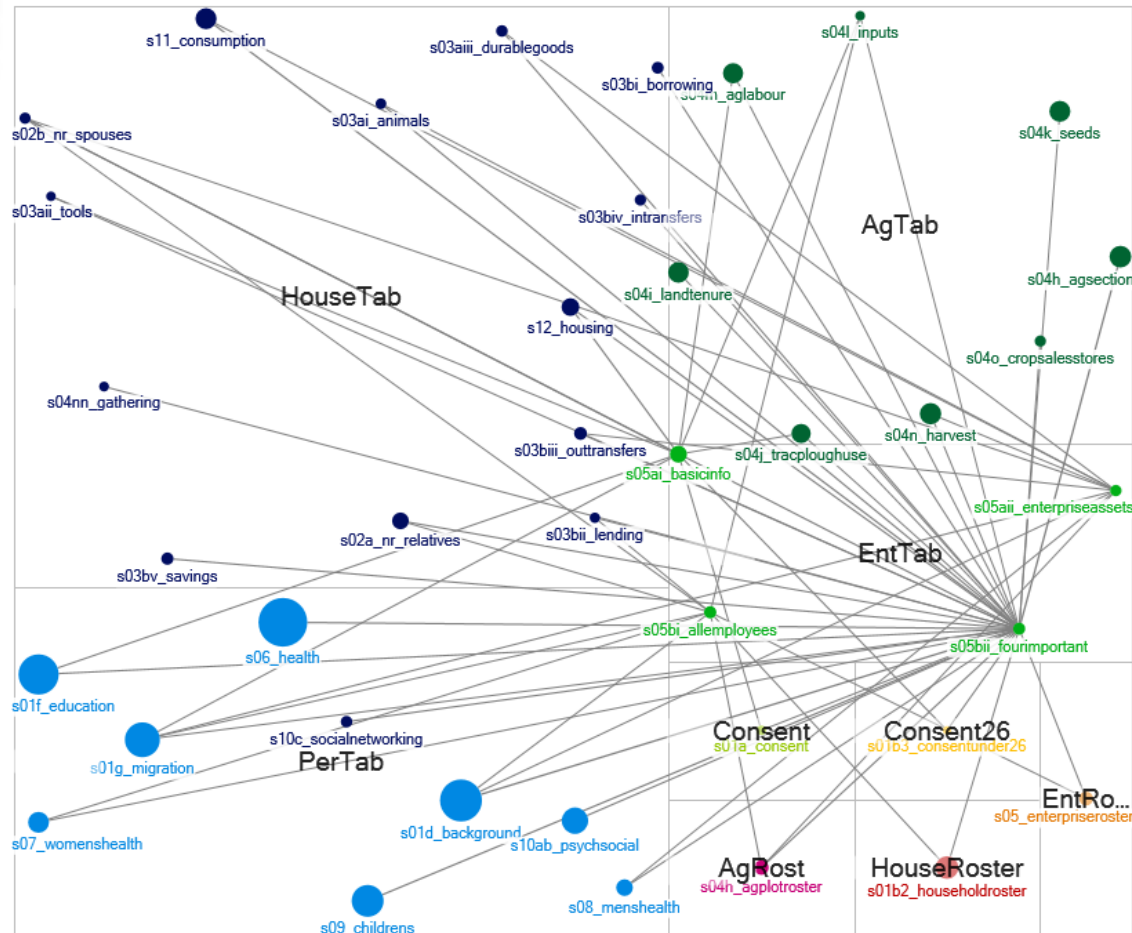
---n/a---

S03Bv: Savings

Done

Most Common Block Moves All Types

- Tendency to work down the columns
- Non-resident Relatives and Consumption introduce multiple common paths

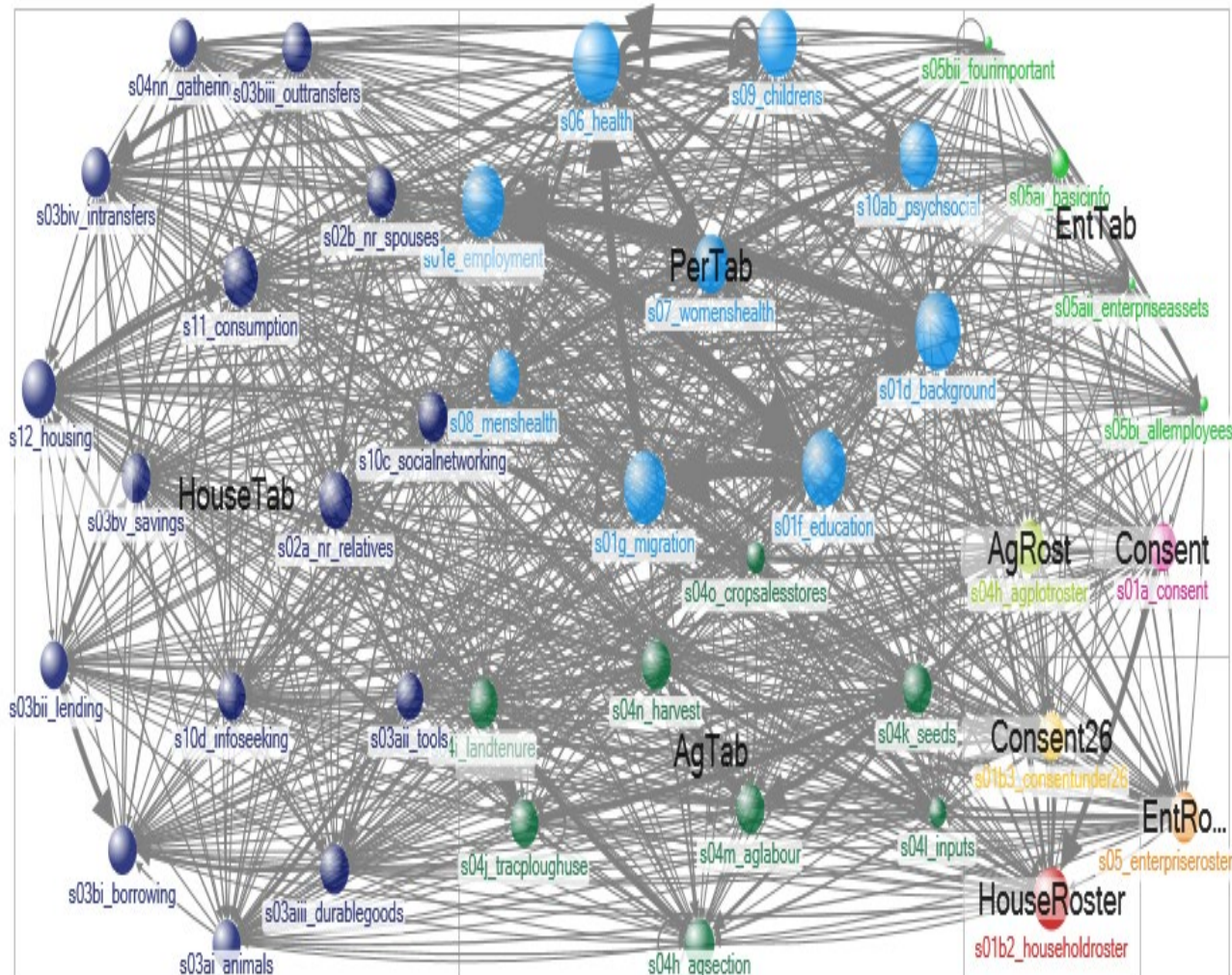


Moving Out of a Section

- interviews showing moves out of the Enterprise section
- “fourImportant” block has most exit moves

Number of Block Moves Per Block

- On a per household basis
 - Average 52 blocks per household
 - Average Block Moves per Block is 1.21
- ☐ Min = 1,
- ☐ Max = 2.82 (for all types)



Instrument Parallel Blocks

- We should have some instructions about the optimal interviewing paths for the desired navigation
 - The parallel blocks programming needs to match with the optimal navigation design
 - The **interviewer training** needs to emphasize the design and avoid “jump around too much”
- How does instrument navigation affect interview length?
 - Order of interview initiation
 - Movements between blocks

Some Observations

- Movement between blocks are with a cost
- Interview length increases with increasing movement between blocks
- Some movements are explainable with the instrument design but others are unsure --- why

Next Steps

- Link the block movement data with interviewer level data to see if there are any connections
- Separate block movement within sections (reasonable) and between sections (why)
- Re-stratify lws by HHs size, R who answered the lw, lw geolocation, or other variable for further analyses
- Have an interviewers debriefing to ask those “why” questions
- Apply all the lessons we learned in this wave to next wave instrument design



The China Health and Retirement Longitudinal Study

- The China Health and Retirement Longitudinal Study (CHARLS) aims to collect a high quality nationally representative sample of Chinese residents ages 45 and older to serve the needs of scientific research on the elderly. The baseline national wave of CHARLS is being fielded in 2011 and includes about 10,000 households and 17,500 individuals in 150 counties/districts and 450 villages/resident committees. The individuals will be followed up every two years. All data will be made public one year after the end of data collection.



Domains of Quality Control (QC) in CHARLES

- Mapping and Listing
- Household Survey
- Biomarkers
- Community Survey

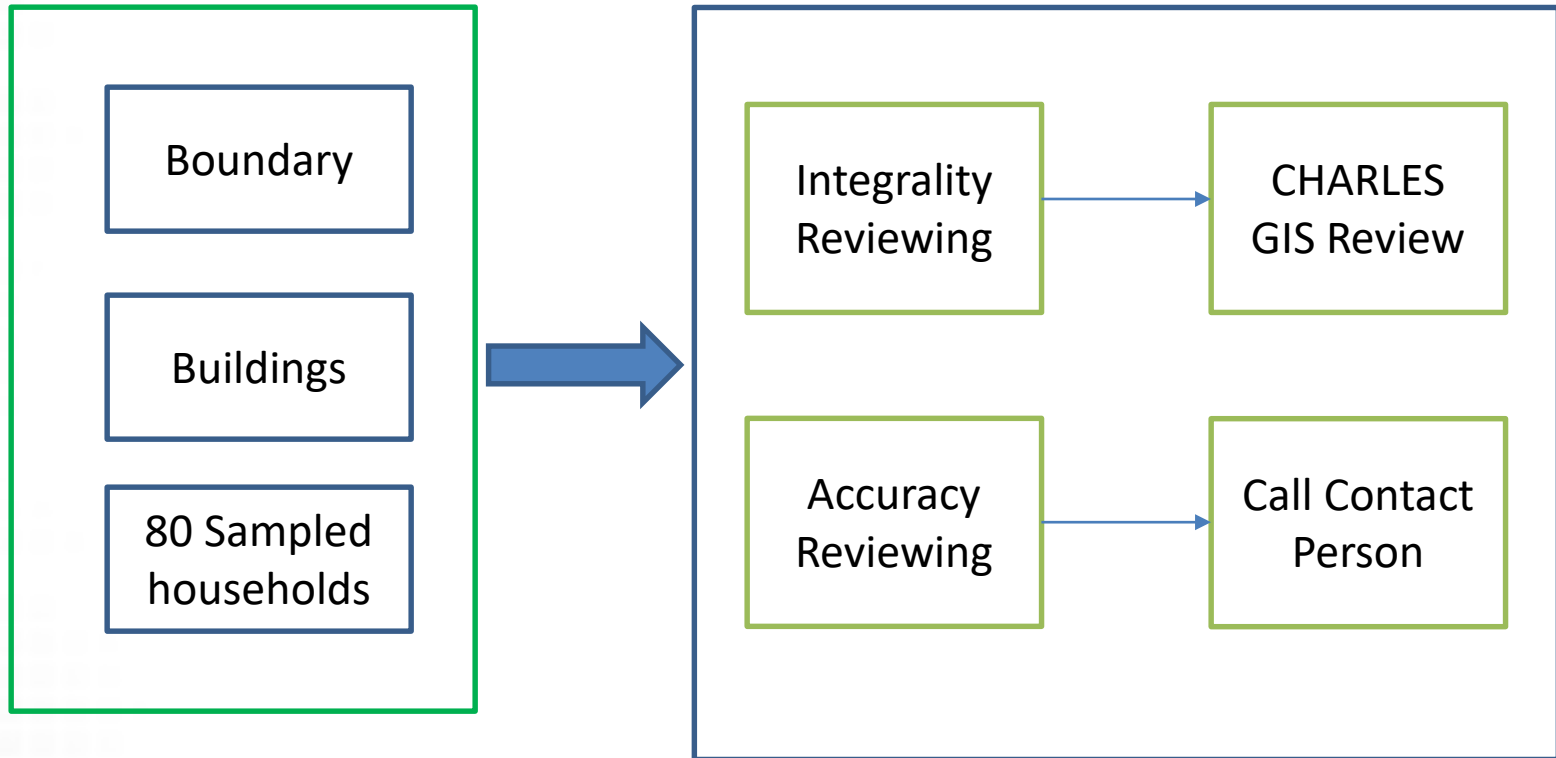


Organizational Structure

- Quality control team
 - Design and implement web-based progress report and QC programming
 - Analyze data
 - Feed information to QC and Field team
- QC supervisors
 - Listen to sound recordings
 - Making telephone calls to respondents
 - Feed information to Quality control team
- Field supervisors
 - Communicate with interviewers; ask for explanations
 - Issue guidelines of work and conduct

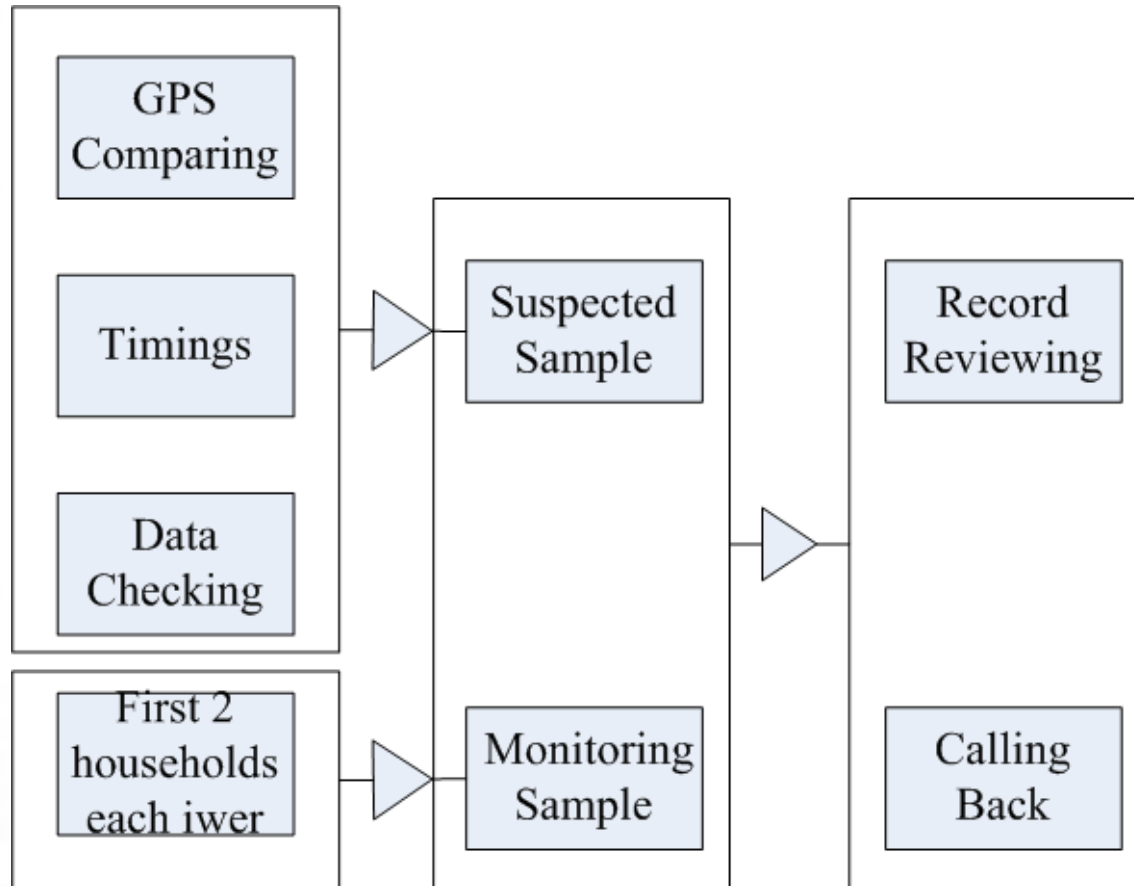


Mapping and Listing





QC in Household Survey





GPS Comparison

- Compare GPS collected in mapping with that collected in household survey
- To guarantee that interviewers went to the correct villages or communities.
- 52.4% GPS successfully collected in household survey



GPS





Module Time

- Pre-determine the minimum time for each module to complete
- Samples falling short of the standard are suspect
- Listen to recording to verify
- Interviewers who do not ask questions in standard ways are warned by supervisors



Check Data

- Checking Points:
 - Key Questions (most important questions)
 - Sensitive Questions (%missing values)
 - Branching Questions (%taking shortcut)
 - Subjective Questions (e.g., mental health)
 - Vignettes (min time)
- Suspect samples are subject to listening recordings and calling



Monitor Samples

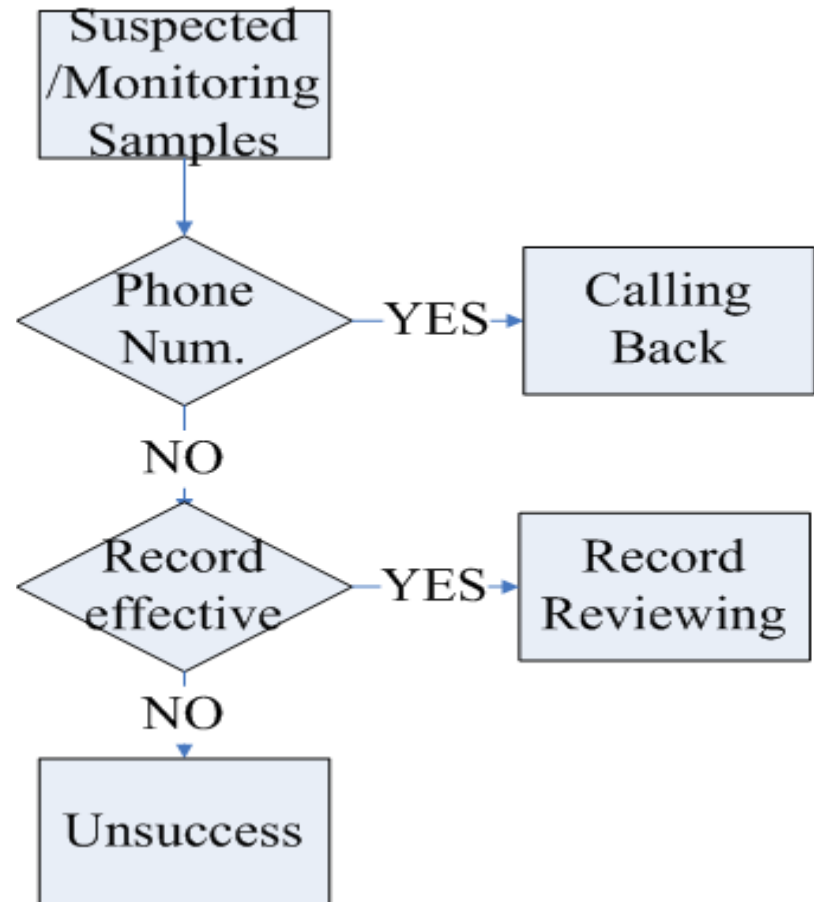
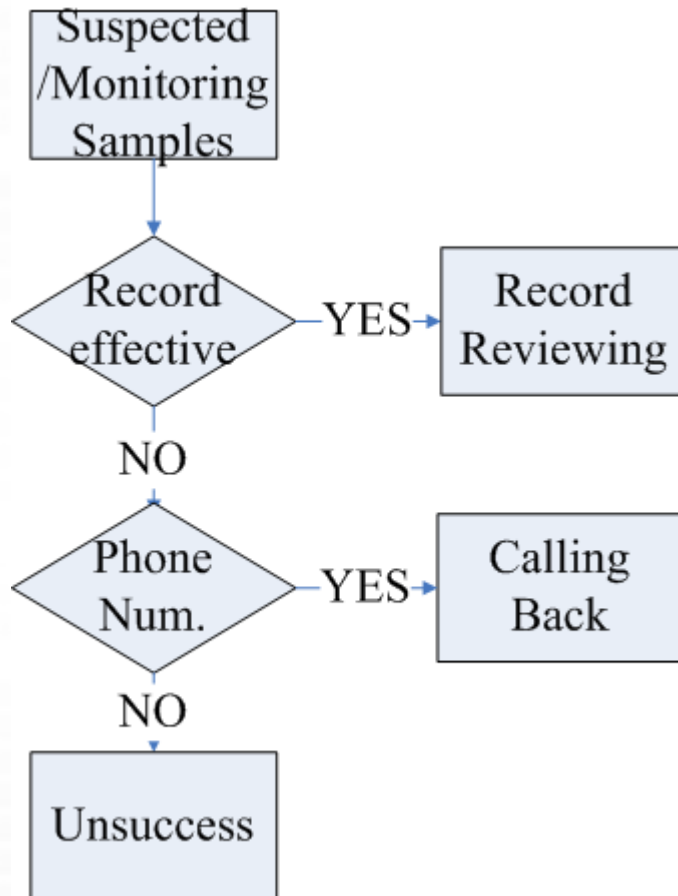
- Samples that are checked even without detecting suspicion
 - Listening to recordings and calling back
- First two completed household surveys of each interviewer
 - to guarantee that each interviewer has samples checked with feedbacks at the beginning of his/her work



Calling Center



Two processes



Interviewers having too many “Unsuccessful” samples are reported to supervisors. There might be equipment problems, refusals of offering contact information or cheating.



What to look for?

- Sound recordings
 - Real Interview
 - Question Jumping
 - Accurate Asking
- Calling back
 - Real Interview
 - Right address
 - Compensations paid or not?
 - A short qxs to re-ask some questions (70% up match is considered ok)

Forthcoming ...

- Photo comparison (doors)



listing (2959)



coverscreen survey



QC in Biomarkers

- Timings: Minimum time requirement of modules such as blood pressure, peak flow test and grip strength
- Sound recording: Introductions at the beginning of each module
- Fingerprints for tracking in 2011

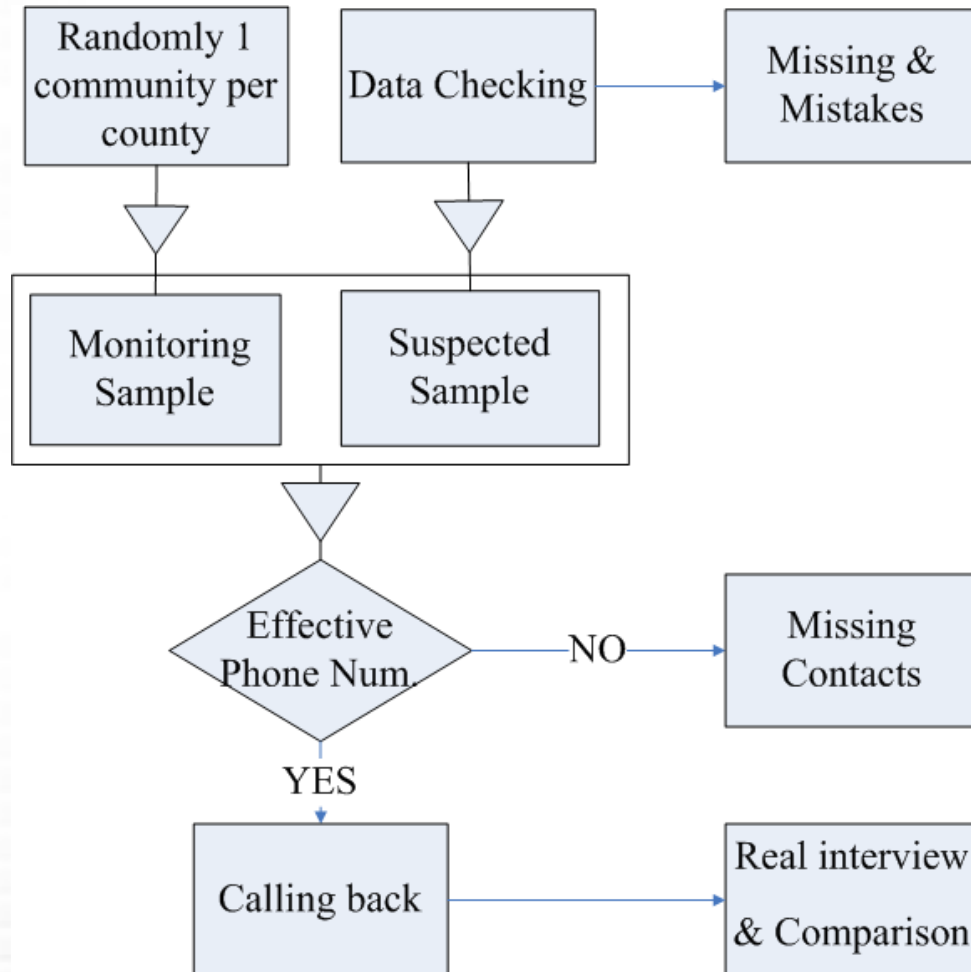


Fingerprint device





QC in Community Survey



Report to supervisors:
samples with too much
missing, mistakes, no
contact information,
unreal interview or
failing in comparison
(below 70%).



Some observations

- Rich paradata collection is becoming the norm
- Paradata can be used across the lifecycle for design issues as well as quality control
- Using paradata for data quality control monitoring is highly effective
- Paradata analysis should be specified throughout the data collection lifecycle but should also have a dynamic component for problem exploration
- Analyzing rich paradata can require a great deal of effort; well designed systems can make a considerable difference



Thank you!

qianyang@umich.edu

bpennell@umich.edu