*24 September 2025*

*13th Meeting of the Expert Group on ICT Household Indicators (EGH), Geneva, Switzerland*

# Real-time quality assurance in ICT household surveys

Dr. Mayank Date, *BDS, MPH*

*Data Scientist, Johns Hopkins Bloomberg School of Public Health, USA*

*Evidence for Digital Transformation (EDiT) Consortium*

School of Public Health
Departement Openbare Gesondheid
Isikolo Sempilo Yoluntu
UNIVERSITY OF CAPE TOWN

JOHNS HOPKINS
BLOOMBERG SCHOOL
*of* PUBLIC HEALTH

*Project funded by*
Gates Foundation

# Data quality assessment processes

**Quality Assurance (QA)**

Processes that ensure adherence to protocols and early error detection.

**Quality Control (QC)**

Techniques used to monitor and verify data

**Quality Improvement (QI)**

Proactive, system-level approach for data-driven efforts that enhance QA and QC systems

# Why do we need high quality data on digital access and use?

- **Paucity of high quality data** on digital access and use, particularly in low resource settings where needed is greatest

- **Data provide important insights into the role technology plays in catalysing health and development**

- **Strategies to improve data quality need to start at survey design and continue through during survey implementation**

# 5 factors that influence survey data quality

**Contextual factors:** Network availability, personnel for supervision, weather, gender dynamics

**Are we asking the right questions?**

**Are the right people prepared to ask the questions?**

**Are the data we are collecting reliable and valid?**

1. Tool Development

2. Enumerator recruitment and training

3. Field Coordination

4. Field Monitoring

5. Data Analytics

**Implementation of AI/ML**

Cognitive Interviews and Pilot testing
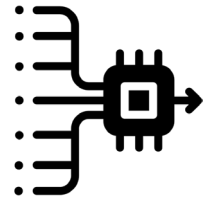
**Data Quality Assessment** → Iteratively feedback results to inform training, supervision, and ongoing implementation

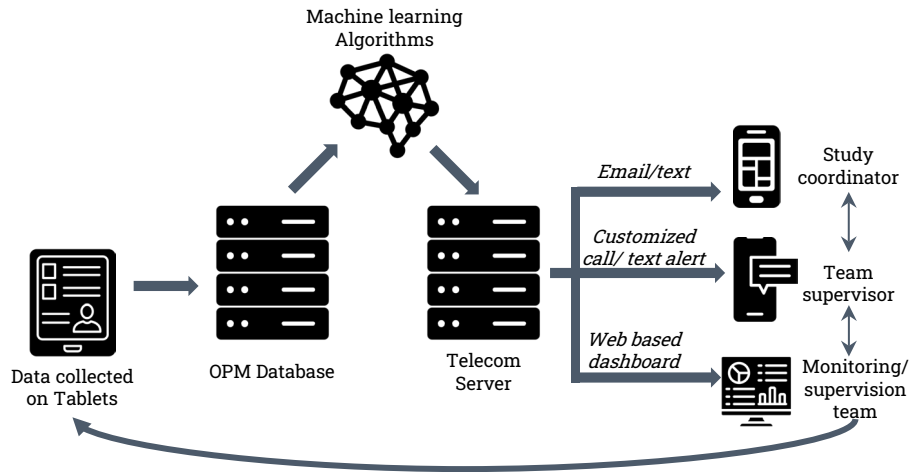# Digitization of data collection presents opportunities to bolster data quality

- Allow for recording and analysis of **meta and paradata** to identify curbstoning and fabrication

- **Machine learning and AI** **can be used for error detection** to improve accuracy and consistency in data
  - Examples of implementation in the pharmaceutical and EHR space
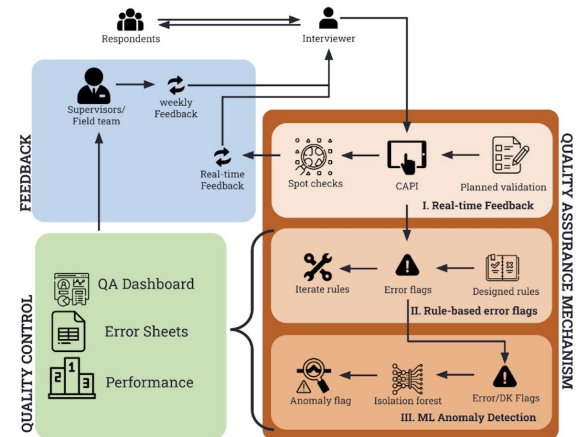  - Limited use in population surveys which represents a missed opportunity

Cohen and Warner - https://www.cambridge.org/core/journals/political-analysis/article/how-to-get-better-survey-data-more-efficiently/6BA27F34CDE007EA8D1AD6692F2E8E60
https://www.sciencedirect.com/science/article/pii/S2667318525000011X
https://www.frontiersin.org/journals/artificial-intelligence/articles/10.3389/frai.2025.1621514/full
https://www.xiahepublishing.com/m/2771-165X/JCTP-2023-00061

Impact evaluation 2018-2021

Digital Access and Use survey in Bihar, India 2024-2025

# Digital access and use survey in Bihar

**Population, 2025**

**Bihar is second most populous state**

with 9.2% of national population [Statistics Times]

**Gross State Domestic Product, 2024-25**

**Bihar ranked 14/36**

at ₹ 9.76 lakh crore ($ 110 billion) [Forbes]

**Human Development Index, 2022**

**Bihar ranked lowest**

of all Indian States [Global Data Lab]

**Gender Inequality Index, 2017-18**

**Bihar ranked 25/36**

of all Indian States [Rural India - Working paper]

# Digital access and use survey in Bihar

**Survey aim:** Measure population level access to and use of mobile phones among men and women 18-60 years of age across 10 districts of Bihar
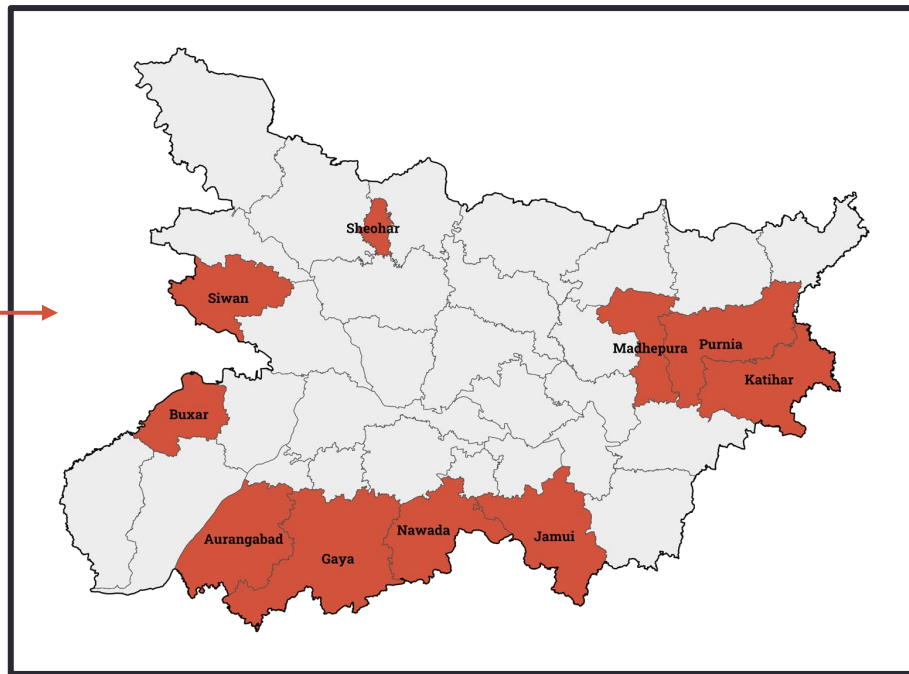
**Data collection:** December 2024-March 2025

**Study sites:** 10 districts in Bihar
- **>68,000 households** listed across 300 villages
- **13,568 respondents** 18-60 years of age including 8567 women and 5,001 men

**Survey Team**
- **57 interviewers** from across India; average of 7.5 years of field experience conducting surveys
- **13 supervisors**
- **5 coordinators**

# Challenges in data collection and QA/QC

- **Poor network and connectivity –** Teams often faced technical challenges like network issues and power outages

- **Low literate population** – concepts assessed not always familiar to respondents

- **Length of survey** – On average, each interview took about **55 minutes** to complete **Target was 3-5 interviews per day**.
  - Poor respondent engagement, especially with no tangible benefits to them
  - Barriers with understanding local dialects and languages
  - Frequent interruptions from family/children
  - Lack of trust

- **Field conditions–**  hot climate, travel distances on foot and by car to reach respondents

- Strict **data storage and sharing policies** prevented implementation of more advanced third-party tools for QA/QC

# Steps for setting up QA/ QC Systems

**Step 1. During tool development, build safeguards into the CAPI system**
- Assess questions for logical skip patterns
- Place time stamps throughout the tool (e.g. start/ stop of sections)
- Assess individual items for logical responses (e.g. age range within 1-100 years)

**Step 2. Develop rule-based error flags and machine learning algorithm protocols**
- Select thresholds and rules for error flags
- Edit rules iteratively during data collection

**Quality Assurance**

# Steps for setting up QA/ QC Systems

**Step 1. During tool development, build safeguards into the CAPI system**
- Assess questions for logical skip patterns
- Place time stamps throughout the tool (e.g. start/ stop of sections)
- Assess individual items for logical responses (e.g. age range within 1-100 years)

**Step 2. Develop rule-based error flags and machine learning algorithm protocols**
- Select thresholds and rules for error flags
- Edit rules iteratively during data collection

**Quality Assurance**

**Step 3. Run data check regularly, track errors and performance**
**Step 4. Generate error reports**

**Quality Control**

# Steps for setting up QA/ QC Systems

**Step 1. During tool development, build safeguards into the CAPI system**
- Assess questions for logical skip patterns
- Place time stamps throughout the tool (e.g. start/ stop of sections)
- Assess individual items for logical responses (e.g. age range within 1-100 years)

**Step 2. Develop rule-based error flags and machine learning algorithm protocols**
- Select thresholds and rules for error flags
- Edit rules iteratively during data collection

**Quality Assurance**

**Step 3. Run data check regularly, track errors and performance**
**Step 4. Generate error reports**

**Quality Control**

**Step 5. Share error sheets with field team for corrective action**

**Feedback**

# QA/QC framework in practice

# I. Real-Time Feedback

Q425 When was the mobile phone within your reach yesterday? In the morning, in the afternoon, in the evening, or in the night?

```
I [Interviewer: Select all that apply. If the phone was
with the respondent for even a part of the time period,
include that time period in the response]
V1 !(self.Contains(6) && self.ContainsAny(1,2,3,4,5))
M1 Not at all can not be select with any other option
V2 !(self.Contains(1) && self.ContainsAny(2,3,4,5))
M2 Whole Day cannot be select with any other option
```
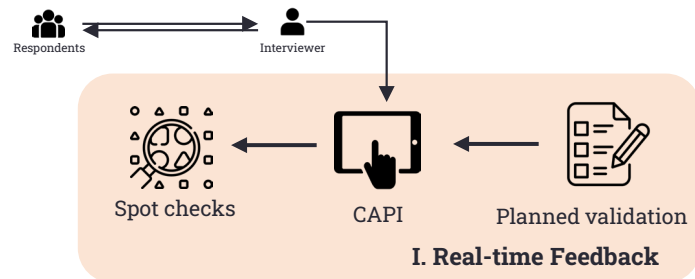
MULTI-SELECT

01 ☐ Whole day
02 ☐ in the morning (6am - 12pm)
03 ☐ in the afternoon (12pm - 6 pm)
04 ☐ in the evening (6pm - 10pm)
05 ☐ in the night (10pm - 6 am)
06 ☐ Not at all



I. Real-time Feedback

- **Validations programmed into the CAPI** software ensuring only valid responses would be accepted
  - Certain age ranges
  - Data format (Numeric/character)
  - Logical selection in multi select questions

- **Spot Checks**
  - About 10%
  - Informed by performance

# II. Rule-based Flagging of Errors - Development

- Rules were defined based on **local context** and internal domain knowledge

- **New rules were iteratively added** as data was assessed weekly



**Respondents** → **Interviewer**

Spot checks ← CAPI ← Planned validation

**I. Real-time Feedback**

Iterate rules ← Error flags ← Designed rules

**II. Rule-based error flags**

| **Contradictory information** | **Implausible response** | **Suspicious practices** |
|---|---|---|
| **Respondent can't read a sentence, but reports having a graduate degree** | **Husband's age is less than 15 or more than 95** | **Interview completed in under 30 minutes** |

# III. Anomaly detection using ML

- Used **Isolation forest** - A type of unsupervised Machine Learning Algorithm to identify anomalous records.
  - **Developed by Liu et al (2008)**
  - **Detecting outliers in the data not caught by rules.**

- **Anomalies** are records that differ from the norm in terms of error patterns or frequent "Don't Know" responses — especially where clear answers are expected.
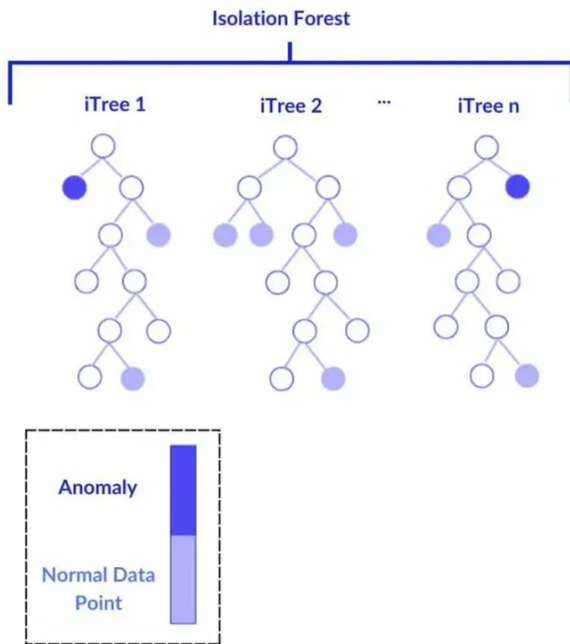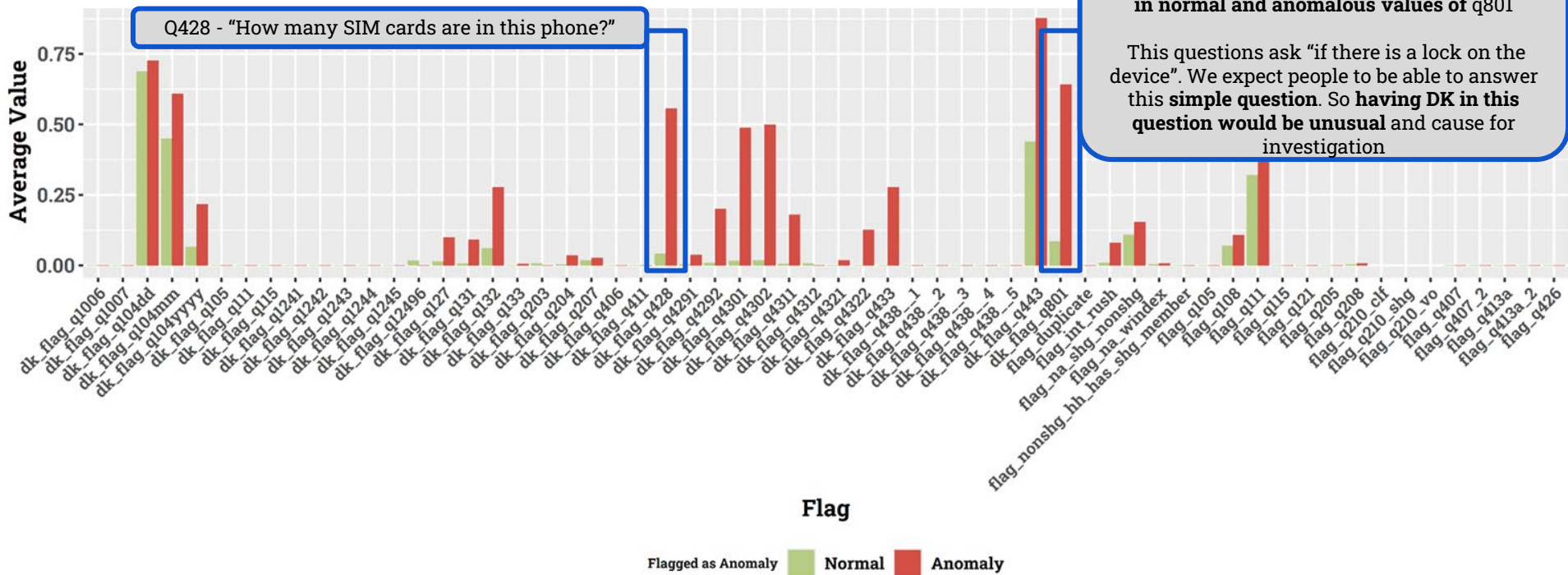


Respondents — Interviewer

**I. Real-time Feedback**
Spot checks — CAPI — Planned validation

**II. Rule-based error flags**
Iterate rules — Error flags — Designed rules

**III. ML Anomaly Detection**
Anomaly flag — Isolation forest — Error/DK Flags

# Strengths of Isolation Forest



Isolation Forest

iTree 1     iTree 2     ...     iTree n

Anomaly

Normal Data Point

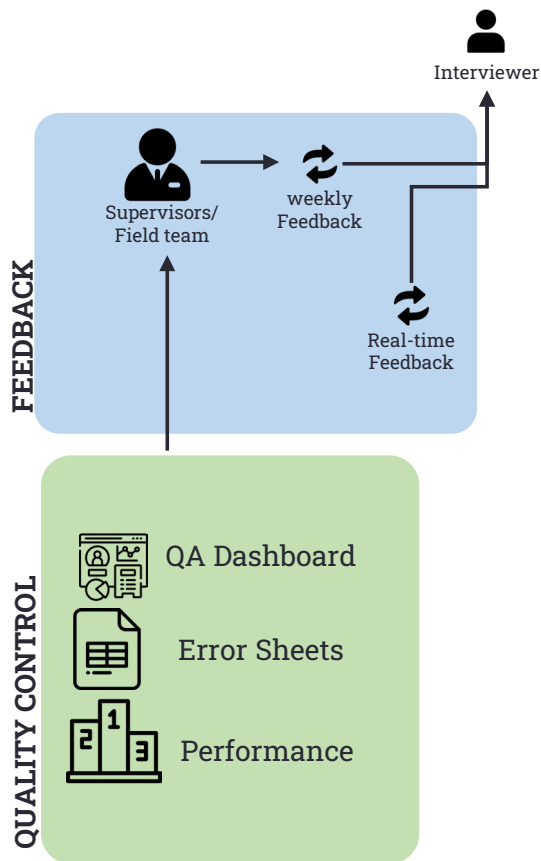Image credit: https://spotintelligence.com/2024/05/21/isolation-forest/

- **Computationally Fast**

- **Generalized.** Does not rely on training data and can accommodate new data easily

- **Low RAM.** Good for large datasets without requiring large computation resource

- Specifically designed for **anomaly detection**

- A 2024 study found it **outperforms\* other models** in detection of fabricated records

*\* Most balanced performance with relatively high precision and recall compared to One-Class Support Vector Machine (SVM), SVM with Stochastic Gradient Descent (SGD), Local Outlier Factor (LOF) algorithm, and Robust Covariance method*

Ref:  https://www.sciencedirect.com/science/article/pii/S2468227624003284

# How the model detects suspicious records

We assessed variables that showed most **difference** in the average values between normal and anomalous records - ie. **Those variables with most unusual discrepancies.**



Q428 - "How many SIM cards are in this phone?"

Eg : The algorithm identified **large difference in normal and anomalous values of** q801

This questions ask "if there is a lock on the device". We expect people to be able to answer this **simple question**. So **having DK in this question would be unusual** and cause for investigation

**Flag**

Flagged as Anomaly    Normal    Anomaly

# Weekly quality control and feedback

Interviewer

**FEEDBACK**

Supervisors/ Field team → weekly Feedback

Real-time Feedback

**QUALITY CONTROL**

QA Dashboard

Error Sheets

Performance

| | | | |
|---|---|---|---|
| **Monday** | **Data download** | Run QA Scripts -Update dashboard -Generate and share error sheets | Review prior week's data quality |
| **Tuesday** | Weekly QA meeting - Dashboard review - Errors discussions - Feedback from field team | | **Field team debrief** |
| **Wednesday** | Update error flags/rules if required based on feedback and analysis | | |
| **Thursday** | **Data download** | Run QA Scripts -Update dashboard -Generate and share error sheets | |
| **Friday** | **Field team debrief** | | |

# Dashboard Set up



## Data Report

Bihar Survey QA Update

PUBLISHED
Saturday Jan 18, 2025 04:16

### Sample Tracking

### i. Progress Tracking

*Tracking completion and response rate of the listing*

**Overall Proportion of SHG Households Listed: 50 %**

| District | HH_Covered (%) | SHG_HH_Listed (%) | SHG_HH_Covered (%) | HH_Listed |
|---|---|---|---|---|
| Buxar | 0 (0%) | 1 (100%) | 0 (0%) | 1 |
| Aurangabad | 601 (9.52%) | 3057 (48.42%) | 302 (9.88%) | 6313 |
| Gaya | 575 (8.51%) | 3408 (50.41%) | 295 (8.66%) | 6760 |
| Katihar | 600 (8.42%) | 3576 (50.17%) | 300 (8.36%) | 7128 |

### ii. Individual Interview Tracking

*No. of valid interviews= ifelse(consent == 1 & interview_status == 100, 1, 0)*

| District | Completed SHG Women | Target SHG Women | Completed Non-SHG Women | Target Non-SHG Women | Completed SHG Men | Target SHG Men | C |
|---|---|---|---|---|---|---|---|
| Buxar | | 375 | | 375 | | 300 | |
| Aurangabad | 403 | 375 | 357 | 375 | 291 | 300 | |
| Gaya | 400 | 375 | 338 | 375 | 263 | 300 | |
| Katihar | 364 | 375 | 347 | 375 | 228 | 300 | |

## Quality Analysis

### a. Cumulative Errors Summary

Search: _____

| enm_name | sup_name | total_cases | total_flags | flag_duplicate | flag_int_rush | f |
|---|---|---|---|---|---|---|
| Amit Goswami | Santanu Maji | 119 | 49 | 0 (0%) | 10 (8.4%) | 0 |
| Prashant Kumar | Jitendra Meena | 159 | 45 | 0 (0%) | 6 (3.8%) | 0 |
| Saumya Patel | Rajesh Kumar Singh | 115 | 45 | 0 (0%) | 8 (7%) | 0 |
| | Sushil | | | 0 (0%) | | |

Showing 1 to 264 of 264 entries

### b. Cumulative Errors Sheet

Search: _____

| enm_name | sup_name | total_flags | flag | link_id | data_source |
|---|---|---|---|---|---|
| Amit Goswami | Santanu Maji | 49 | flag_int_rush | 2_18_241_162_1_0 | wom_data |
| Amit | Santanu | 49 | flag_q108 | 2_18_241_162_1_0 | wom_data |

### c. Weekly Error Summary

#### Week: 7 (18 Jan - 24 Jan)

Total flags this Week: 79

**Flag Counts**

Search: _____

| flag_name | count |
|---|---|
| flag_anomaly | 31 |
| flag_na_windex | 15 |
| flag_q407 | 12 |
| flag_q413a | 7 |
| flag_q208 | 3 |

Showing 1 to 10 of 10 entries

**Weekly Errors Summary for each Enumerator**

Search: _____

| enm_name | sup_name | total_cases | total_flags | flag_duplicate | flag_int_rush | f |
|---|---|---|---|---|---|---|
| Ghanshyam Sahu | Bibekananda Prusti | 4 | 3 | 0 (0%) | 0 (0%) | 0 |
| Gulshan Kumar | Sushil Kumar | 23 | 3 | 0 (0%) | 0 (0%) | 0 |

### d. Enumerator Performance



Survey-Wide Issues Over Time

Proportion of Error cases each week - Excluding Supervisors and under 5 total cases



Total Interview Time Distribution

```
clip_dis_time_box = dis_time_box +
    coord_cartesian(ylim = c(0, 120))
clip_dis_time_box
```

Distribution of Interview Times by District and SHG HH Status

# Weekly Error Sheets

- Field team used this to investigate and correct errors on a case by case basis.
- Example: Week 5 Error Sheet

# 85% improvement in error rates over 12 weeks



Weekly Error Trends by Data Source and Flag

2 new flags introduced - contributes to spike in error count across all data sets

Introduced a flag for respondent with higher education being reported poorly literate.

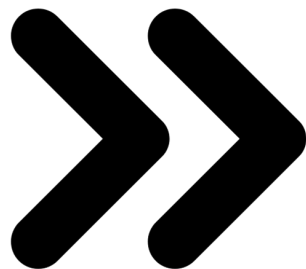**This seemed to disproportionately affect women's data**

# Key Performance Indicators

- All sample size and data collection targets were met within scheduled time frame

- Observed **85% drop in error rates** across the course of data collection

- Observed **20 - 50 % point decrease in error rates** for interviewers with higher error rates at start of data collection
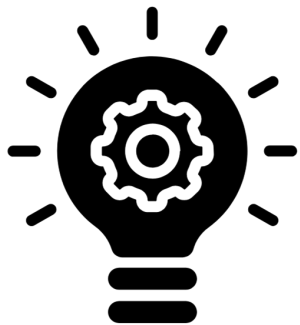
# Next steps

- **Integrate paradata -** Keystroke tracking, timestamps, GPS tracking, and audio capture **where possible**.

- **AI/LLM integration**. Explore use of large language models for real-time curbstoning detection, ensuring privacy safeguards.

- **Open-source resources.** Share rules and anomaly detection pipelines to enable replication and scaling.

# Anticipated use for QA / QC approach

- **Improve survey data quality**

- **Opportunity to apply this approach to improve data quality more widely to a range of development programs**
  - End to end digitization of health systems
  - Government efforts in India to digitize Self-help groups for economic empowerment

# Thank you
## on behalf of the EDiT Team

**Dr. Mayank Date**

**mdate1@jhu.edu**

*Data Scientist, Johns Hopkins University*

Project funded by
## Gates Foundation

*evidence-digital.org*

Scan to learn more about
**Evidence for Digital Transformation (EDiT)**