# Mobile SMS survey data management and preservation

Inna Kouper

Charitha Madurangi, Kunalan Ratharanjan, Tom Evans, Beth Plale
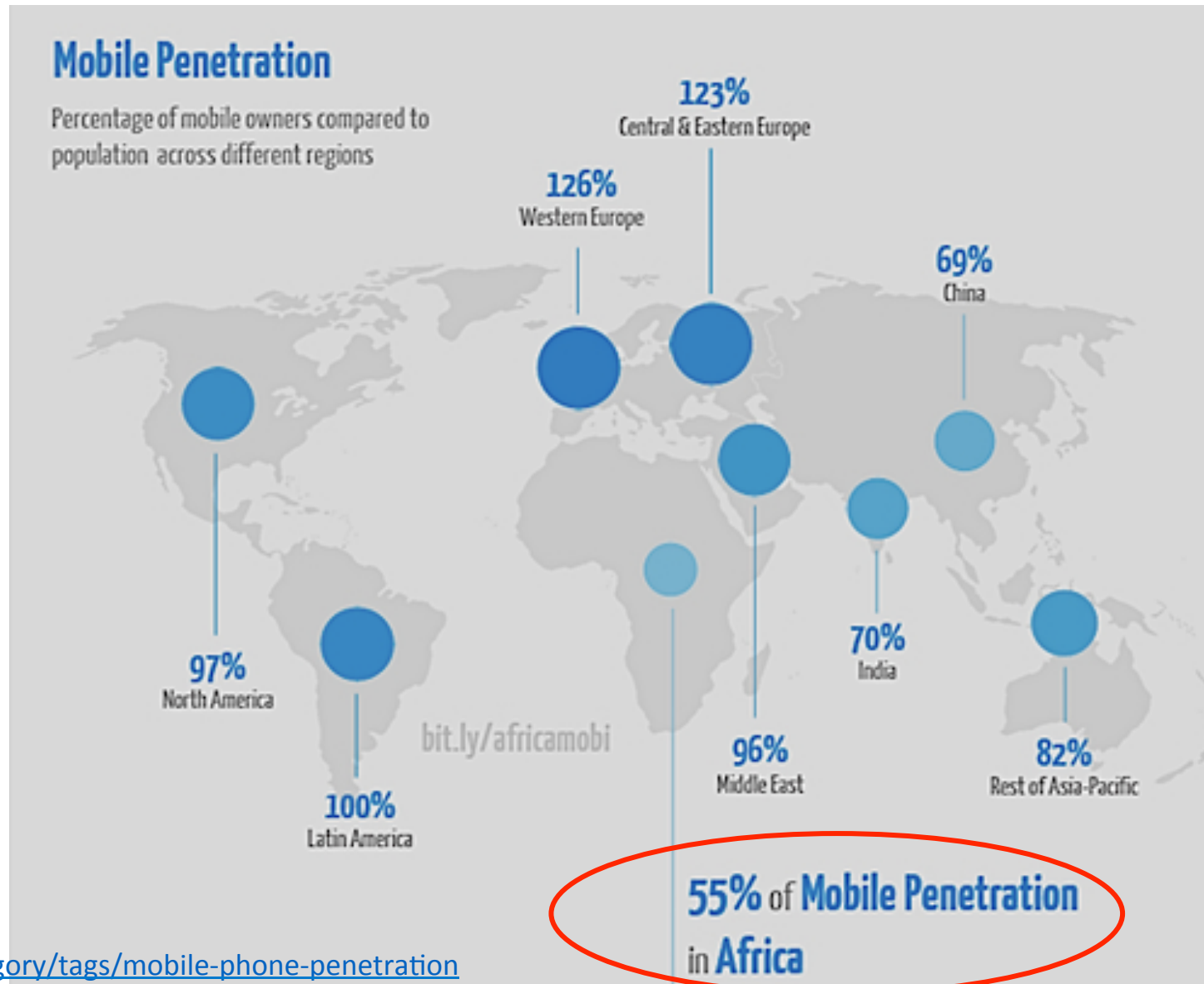
Indiana University

National Science Foundation
WHERE DISCOVERIES BEGIN

# Mobile phone penetration



**Mobile Penetration**

Percentage of mobile owners compared to population across different regions

- 123% Central & Eastern Europe
- 126% Western Europe
- 69% China
- 97% North America
- 100% Latin America
- 96% Middle East
- 70% India
- 82% Rest of Asia-Pacific
- 55% of Mobile Penetration in Africa

bit.ly/africamobi

# Mobile phones in Africa

## Cell Phone Ownership Surges in Africa

*Adults who own a cell phone*



U.S. 89
S. Africa 89
Ghana 83
Kenya 82
Tanzania 73
Uganda 65

100%

64

50

33

10
10
9
8
0

2002    2007    2014

Note: U.S. data from Pew Research Center surveys.
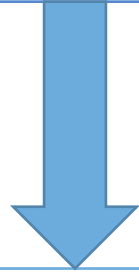
Source: Spring 2014 Global Attitudes survey. Q68.

**PEW RESEARCH CENTER**

http://www.pewglobal.org/2015/04/15/cell-phones-in-africa-communication-lifeline/africa-phones-7/

# Project: Agricultural Decision Making and Food Security in Africa

The project examines how small-scale farmers adapt to food and climate variability.
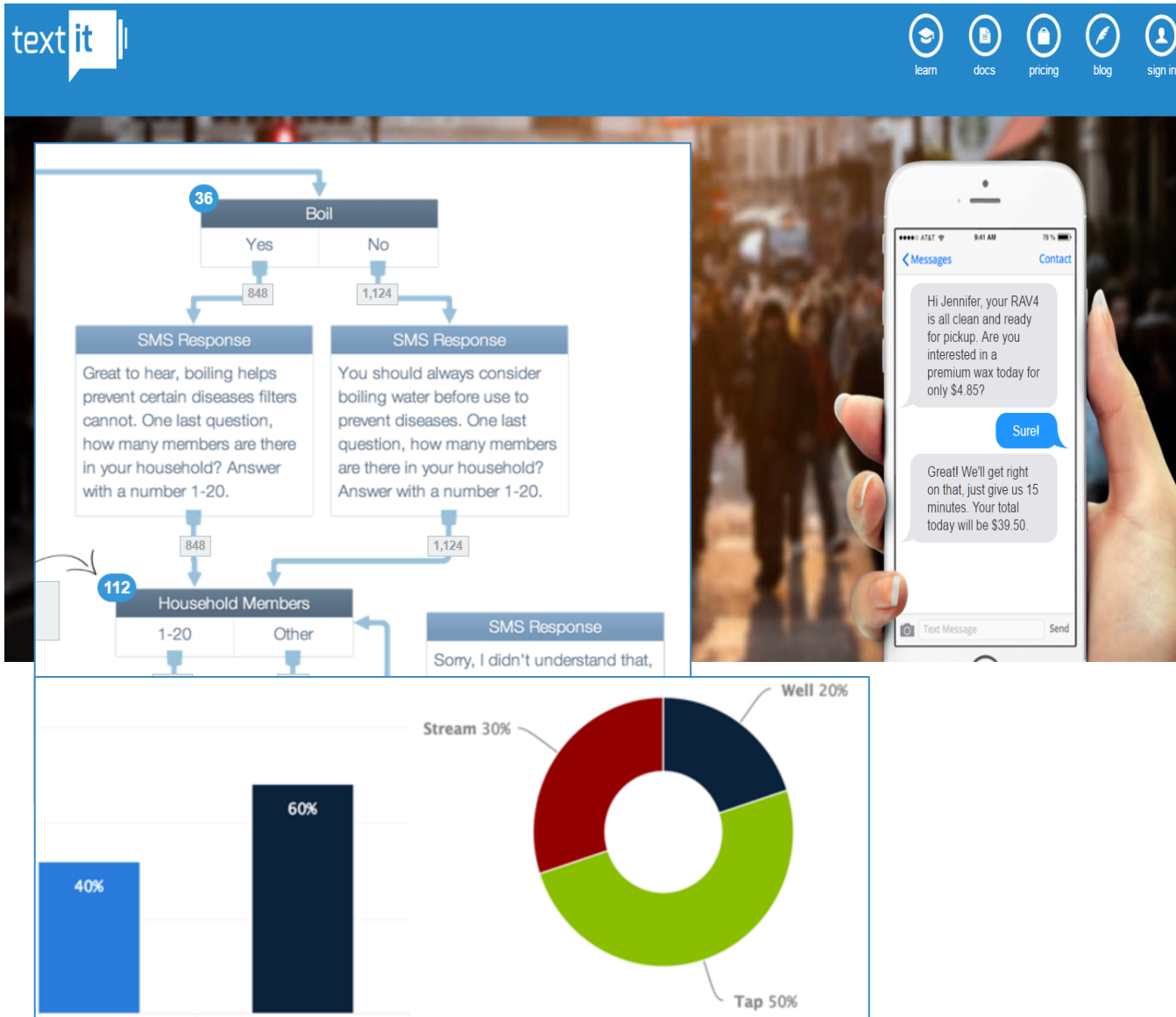
It integrates physical models with real-time environmental data and weekly farmer decision making in individual fields.

Farmers are asked weekly about their decisions to plant, grow, and harvest and about weather conditions.

# TextIt SMS Platform



- Cloud-based commercial SMS service
- Builds surveys through GUI
- Provides summaries and minimal analytics
- Data can be downloaded manually or via API

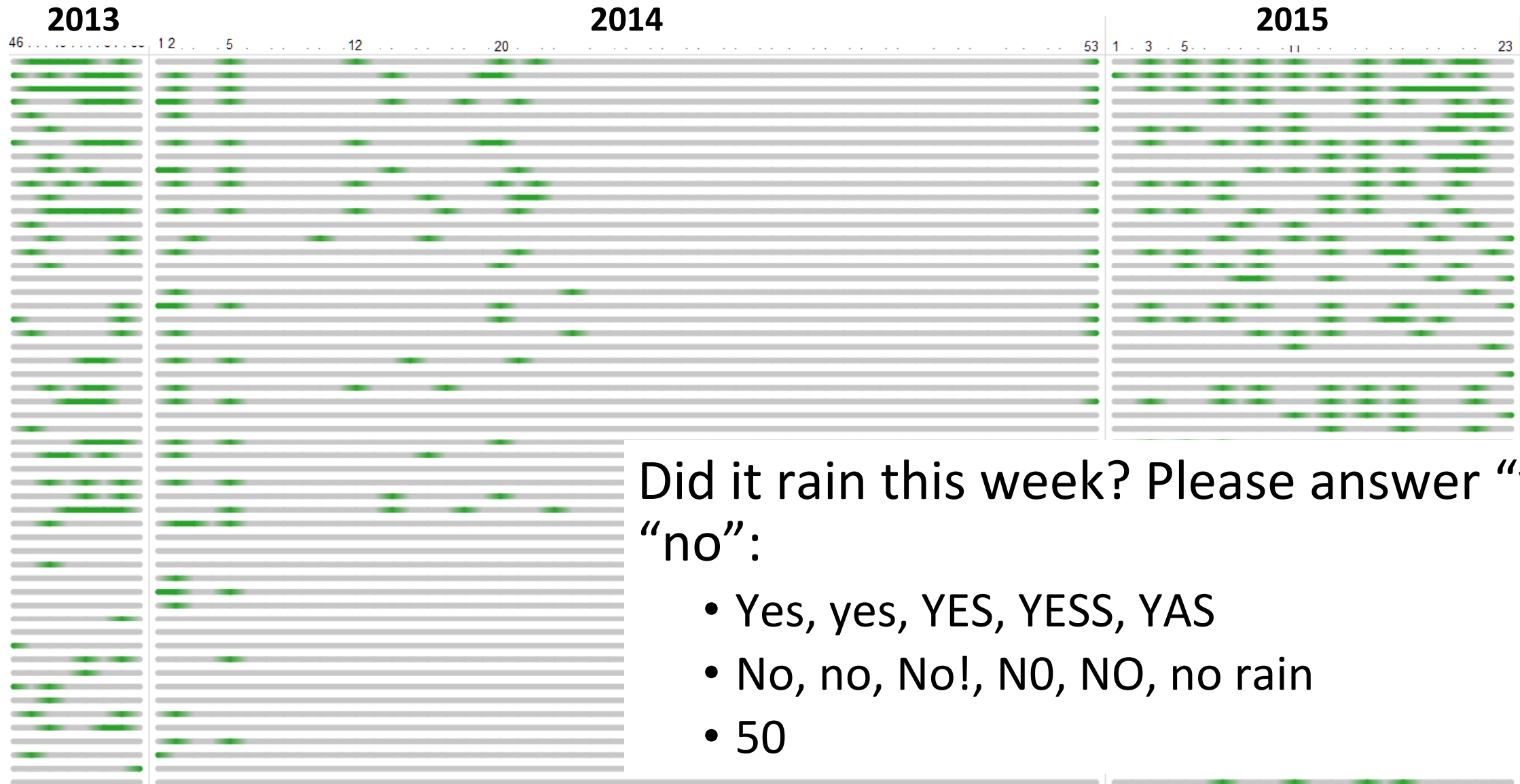# Text messaging (SMS) to collect data

## Pros

- High-frequency, automated data collection

- Large sample sizes

- Relatively low cost

## Challenges

- Set up learning curve

- Data may be incomplete due to non-response or lack of SMS credits

- Typing increases errors

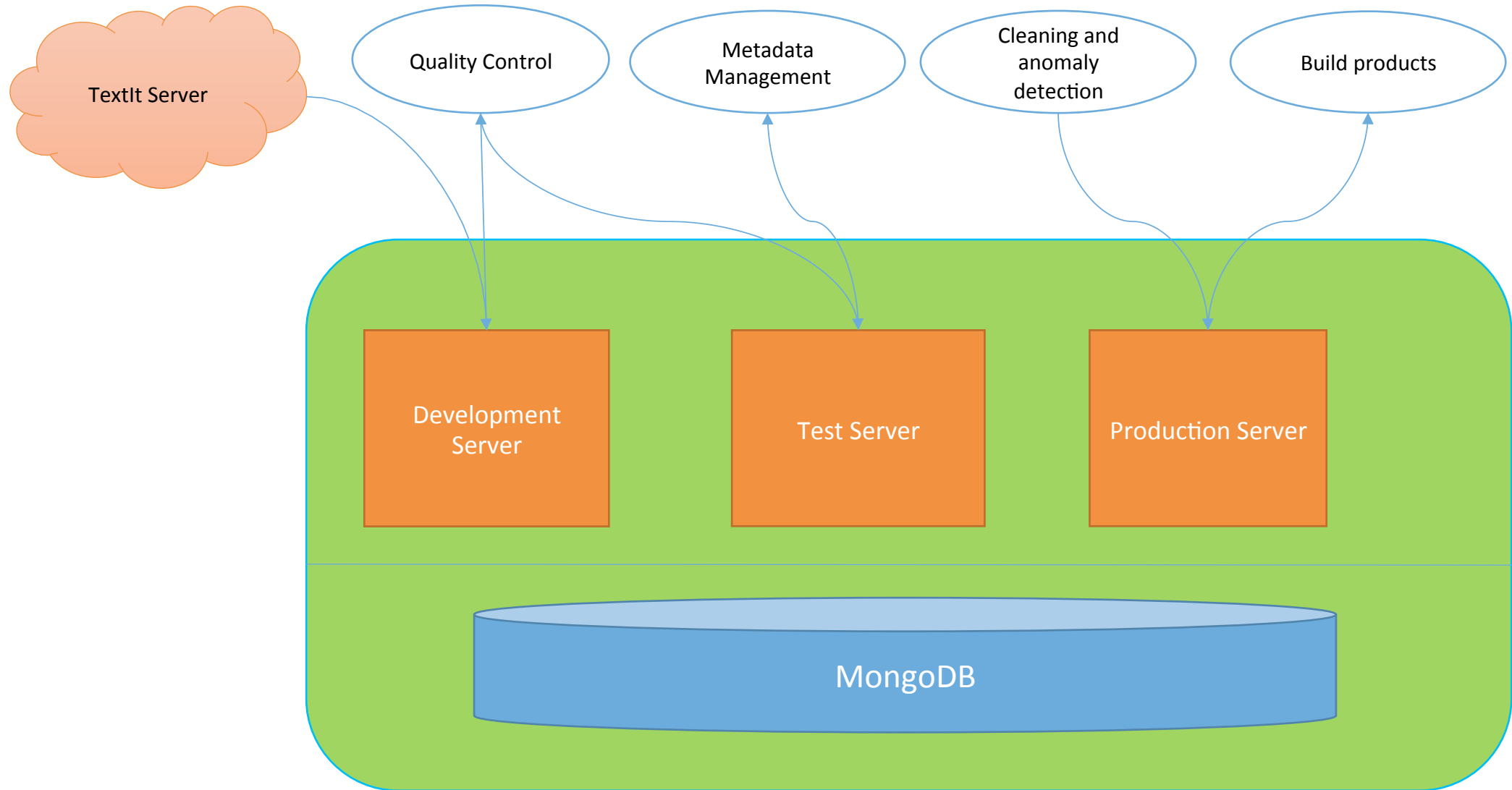- Some programming expertise required

- Possibilities depend on platform

# Data completeness and errors

**2013**         **2014**         **2015**



Did it rain this week? Please answer "yes" or "no":

- Yes, yes, YES, YESS, YAS
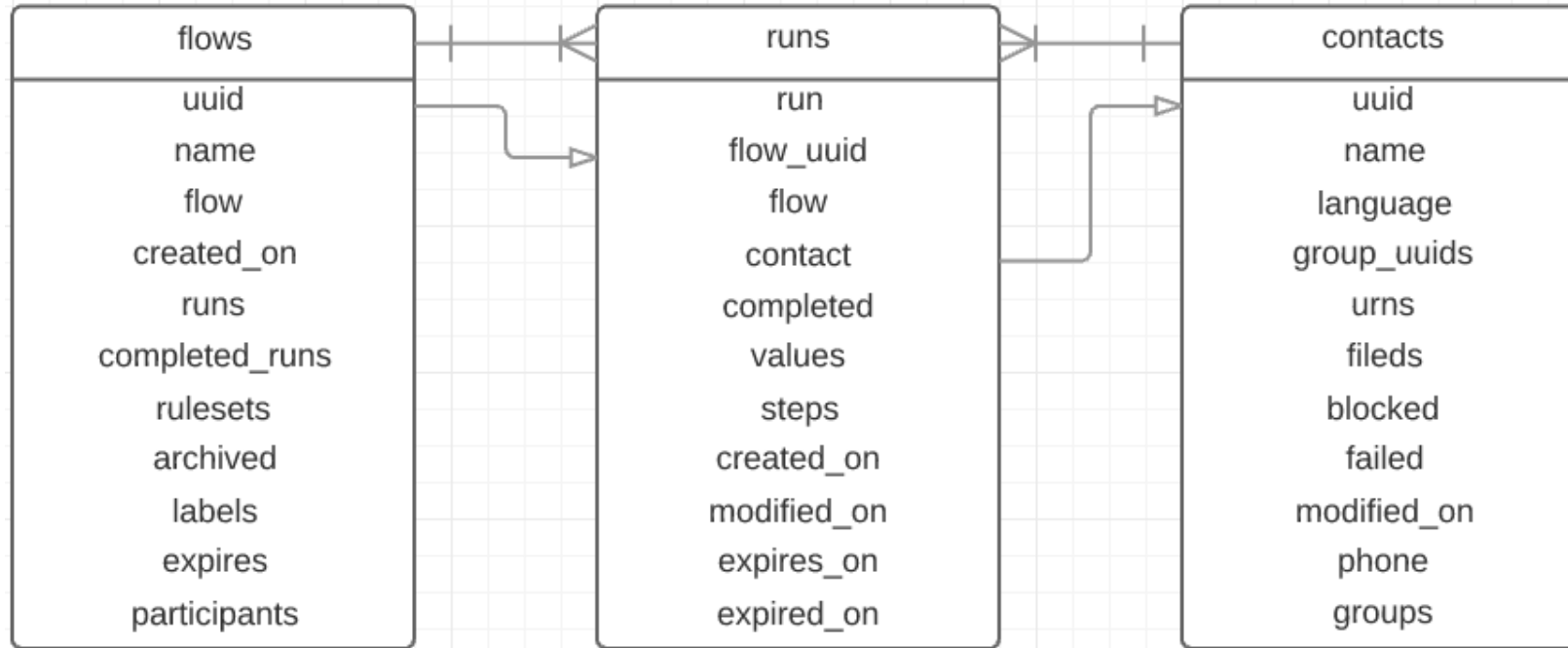- No, no, No!, N0, NO, no rain
- 50

# Data retrieval via API

```
"count": 23,
"next": null,
"previous": null,
"results": [{
    "flow_uuid": "0f6ded36-0ed4-4c10-86f5-eaa6b0fd0f9c",
    "flow": 29541,
    "run": 2110573,
    "contact": "7ed0ad89-969d-4cce-b012-c7d91d9ae731",
    "completed": true,
    "values": [{
        "category": {
            "base": "No"
        },
        "node": "b7f97b23-3c3c-4baa-b34c-11d997dba5e2",
        "time": "2015-04-29T10:09:49.894Z",
        "text": "NO",
        "rule_value": "NO",
        "value": "NO",
        "label": "mmf harvest"
```

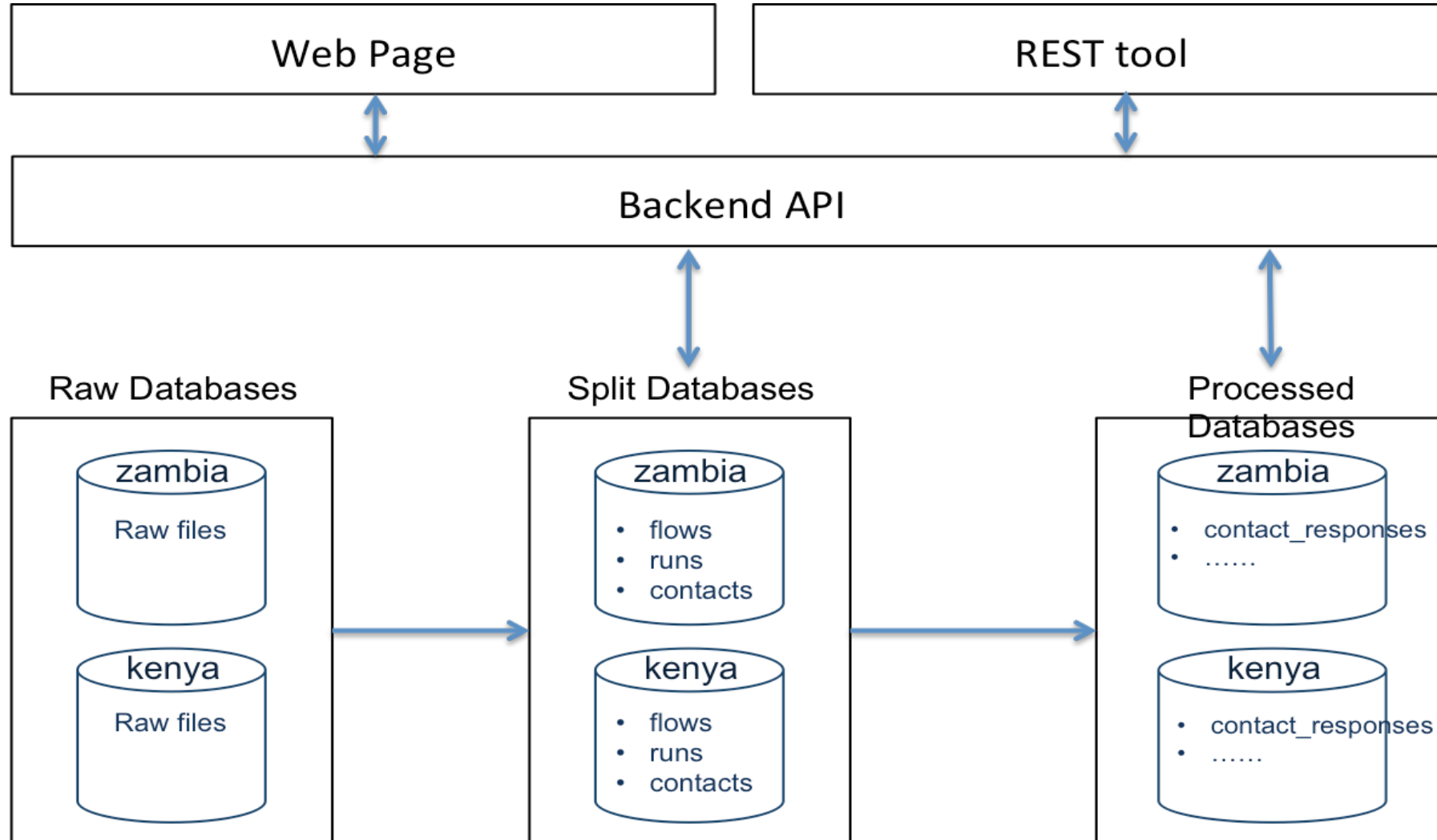# Our Approach: Data Management and Preservation Pipeline
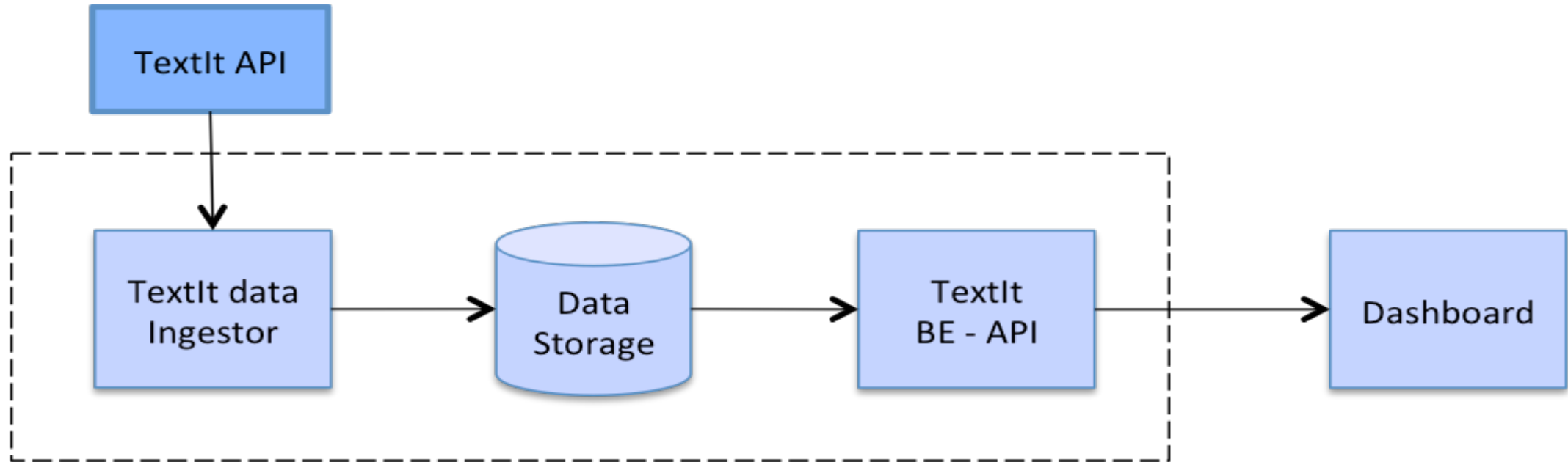
# Database Design



Data is organized into three collections – *Raw*, *Split*, and *Pre-processed* for two countries

- *Raw* preserves original data
- *Split* organizes data into flows (surveys), runs (responses within one survey) and contacts (respondents)
- *Pre-processed* minimizes workload for queries

# Storage and Access Architecture

# Data Retrieval Pipeline



- Pipeline is fully automated
- Data from TextIt is ingested weekly
- Metadata from TextIt is stored in addition to data (survey responses)
- Data is restricted to users within university
- Data can be downloaded via dashboard or http request

# Metadata Improvements

**Retrieved from TextIt**

- IDs
- Flow names
- Total number of runs (responses)
- Number of completed runs
- Variable labels

**Added by the team**

- Country
- Season
- Creator
- Date created
- Run start and end date / time
- Flow type
- List of questions
- Farmer contact and location

# Data / Metadata GUI



- View and modify metadata details for flows and contacts
- Download inactive contacts for a given time period
- Download farmer response details for given question and time period

# Monitoring Dashboard



- Flows details
- Farmer contacts details
- Flow completion rates
- Inactive / non-responsive contacts

# What have we achieved?

Curation, curation, curation!

- **Automated pipeline** enables consistent long-term preservation of raw data and full control over it.

- **Pre-processing and transformations** help to quickly retrieve subsets of data for analysis.

- **Improved metadata** facilitates search, access, and future re-use.

- **Organized storage** enables future visualizations and integrations.

- **Interface and dashboard** makes interactions with data easier, no technical skill required

# Challenges

Curation, curation, curation!

- **Security** – need to add proper authentication and access options for research team and the public

- **Data cleaning** – some can be automated, but most is still manual

- **Maintenance**
  - Changes in commercial platforms (e.g., APIs) require modifications of backend and data
  - Staff is needed for ongoing curation and technical maintenance
  - Data preservation is not central to research projects

- **Analysis and dissemination** – still done outside of the pipeline, not reproducible

# Future Work (Questions)

- What standards in data documentation and preservation can help to improve this work?
- How can SMS data be integrated with other types of data (e.g., sensor data or household interviews)?
- What analytical products are most useful and for what types of stakeholders?
- **How can we measure the impact from curation?**