# The Rakai Health Sciences Program (RHSP):
## Developing *a Data Warehouse*



By: ANTHONY NDYANABO
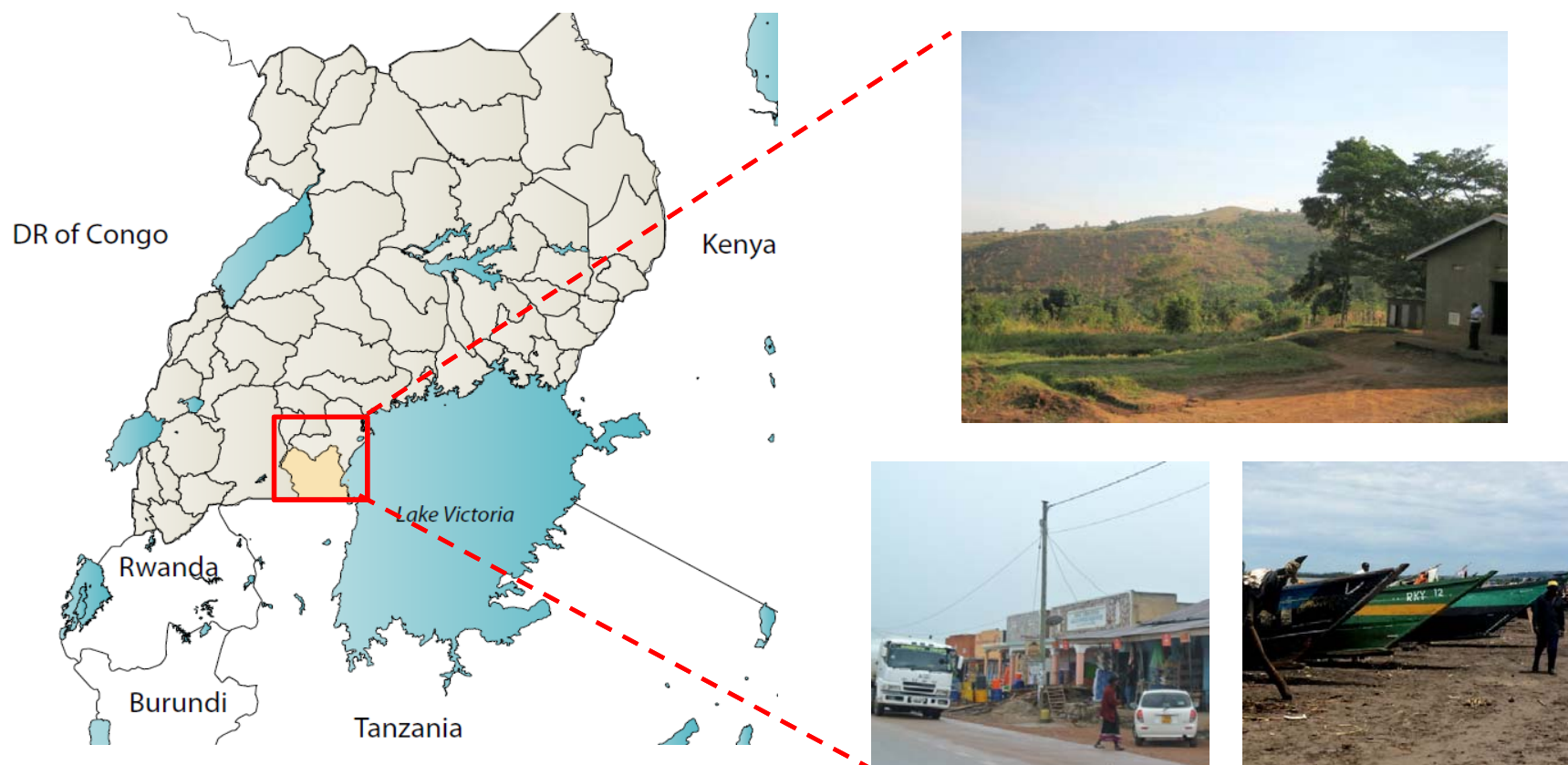
# The Rakai Health Sciences Program (RHSP)

- The RHSP is currently one of the largest and oldest community-based non-profit research endeavors on HIV/AIDS, associated infections and reproductive health, in Africa.

- The program is currently conducting HIV prevention studies, basic and clinical laboratory research, and operations research/evaluation to develop improved strategies of service delivery in southern Uganda.

# Rakai Region, Uganda

- Agrarian, trading, and fishing communities.

# Why Rakai?

- D Serwadda, N Sewankambo, *et al* identified "slim disease" in Rakai in 1982.

**SLIM DISEASE: A NEW DISEASE IN UGANDA AND ITS ASSOCIATION WITH HTLV-III INFECTION**

| | |
|---|---|
| D. SERWADDA | R. D. MUGERWA |
| N. K. SEWANKAMBO | A. LWEGABA |
| J. W. CARSWELL | G. B. KIRYA |
| A. C. BAYLEY | R. G. DOWNING |
| R. S. TEDDER | S. A. CLAYDEN |
| R. A. WEISS | A. G. DALGLEISH |

*"A new disease has recently been recognized in rural Uganda."*

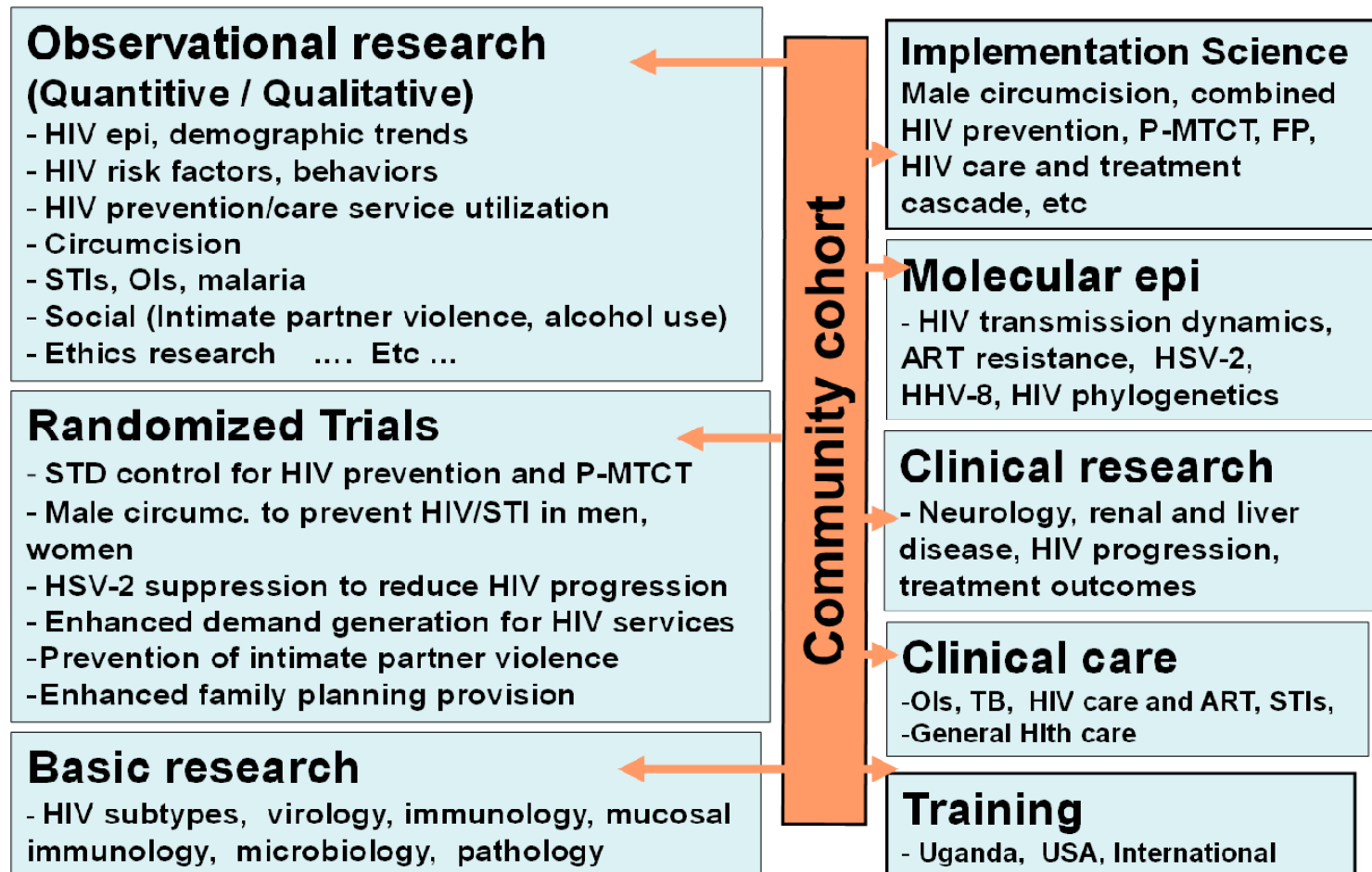Serwadda D, Sewankambo N, et al. Lancet 1985.
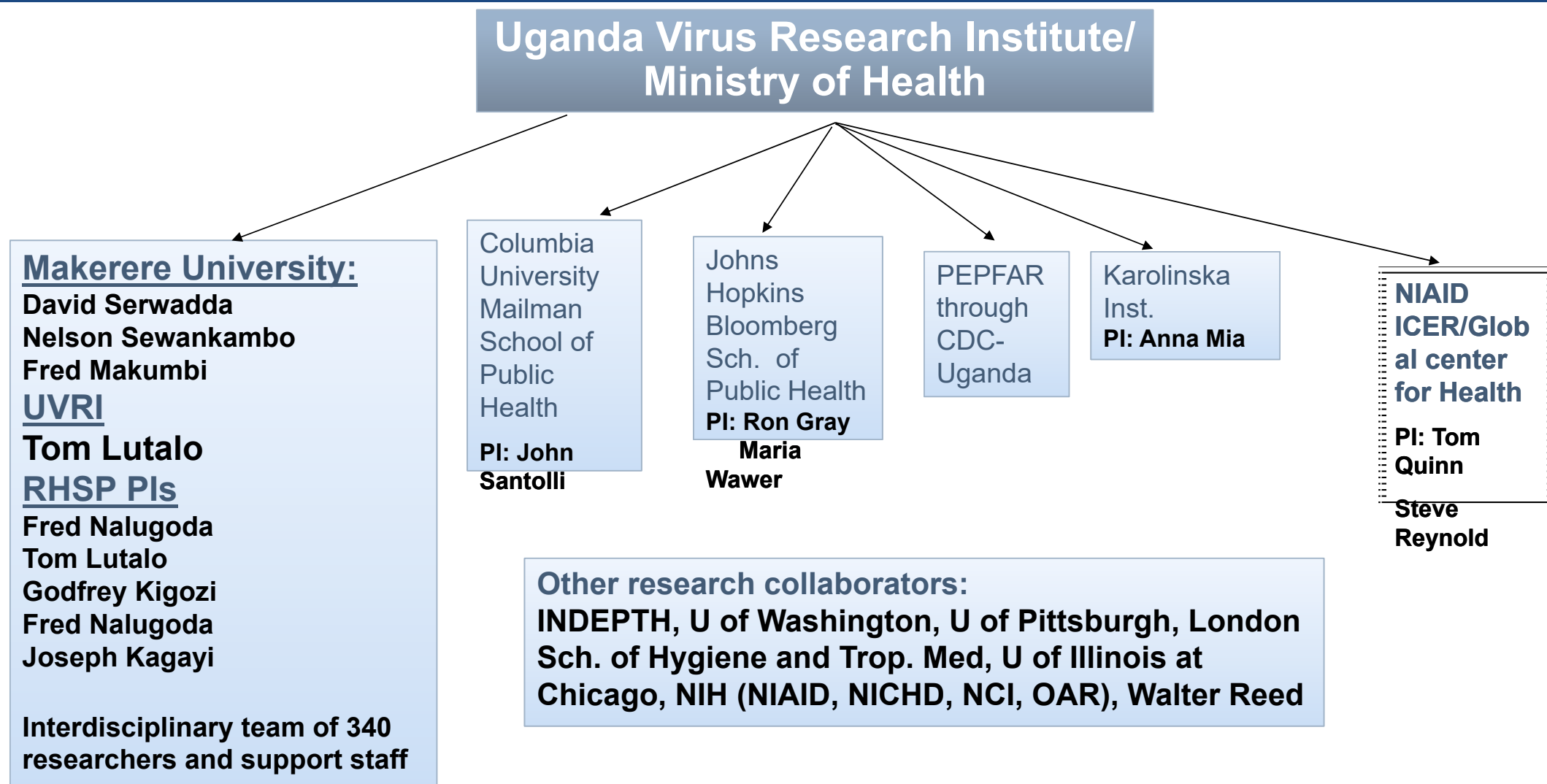
# The Rakai Community Cohort Study (RCCS)

- Open population-based **household census and cohort survey** of 40 communities ongoing since 1994 conducted by the RHSP

- **Household Census**
  - All household residents included irrespective of age
  - <u>Data obtained:</u> residence status, births and deaths, family relationships, household assets/wealth, GPS coordinates for households and local built environment (e.g. schools, bars, hotels).



- **RCCS Cohort Survey**
  - Currently restricted to ages 15-49 years
  - <u>Questionnaire</u>: Detailed demographic, sexual network and behaviors, health status, and health care utilization questions
  - <u>Specimen collection</u>: Serum/plasma for HIV, STI, and other tests; buffy coat for human genetic testing

# Rakai Community Cohort Study (RCCS)

**Observational research**

**(Quantitive / Qualitative)**
- HIV epi, demographic trends
- HIV risk factors, behaviors
- HIV prevention/care service utilization
- Circumcision
- STIs, OIs, malaria
- Social (Intimate partner violence, alcohol use)
- Ethics research    .... Etc ...

**Randomized Trials**
- STD control for HIV prevention and P-MTCT
- Male circumc. to prevent HIV/STI in men, women
- HSV-2 suppression to reduce HIV progression
- Enhanced demand generation for HIV services
- Prevention of intimate partner violence
- Enhanced family planning provision

**Basic research**
- HIV subtypes, virology, immunology, mucosal immunology, microbiology, pathology

**Community cohort**

**Implementation Science**
Male circumcision, combined HIV prevention, P-MTCT, FP, HIV care and treatment cascade, etc

**Molecular epi**
- HIV transmission dynamics, ART resistance, HSV-2, HHV-8, HIV phylogenetics

**Clinical research**
- Neurology, renal and liver disease, HIV progression, treatment outcomes

**Clinical care**
-OIs, TB, HIV care and ART, STIs,
-General Hlth care

**Training**
- Uganda, USA, International

# RHSP: A global collaboration

**Uganda Virus Research Institute/ Ministry of Health**

**Makerere University:**
**David Serwadda**
**Nelson Sewankambo**
**Fred Makumbi**
**UVRI**
**Tom Lutalo**
**RHSP PIs**
**Fred Nalugoda**
**Tom Lutalo**
**Godfrey Kigozi**
**Fred Nalugoda**
**Joseph Kagayi**

**Interdisciplinary team of 340 researchers and support staff**

Columbia University Mailman School of Public Health

**PI: John Santolli**

Johns Hopkins Bloomberg Sch. of Public Health
**PI: Ron Gray**
**Maria Wawer**

PEPFAR through CDC-Uganda

Karolinska Inst.
**PI: Anna Mia**

**NIAID ICER/Global center for Health**

**PI: Tom Quinn**

**Steve Reynold**

**Other research collaborators:**
**INDEPTH, U of Washington, U of Pittsburgh, London Sch. of Hygiene and Trop. Med, U of Illinois at Chicago, NIH (NIAID, NICHD, NCI, OAR), Walter Reed**

# RCCS Data

- Data collected thus far**:**
  - 18 completed survey/census rounds (19th ongoing)
  - ~22,000 study participants currently
  - ~80,000 total participants
  - >500K archived laboratory samples

- Historically, data for each survey round was stored in own separate FoxPro Database (hundreds of tables!)

- Data requests can be overwhelming
  - Assembling longitudinal data sets challenging!

# Why the data warehouse?

- To consolidate the available cohort study data rounds into a single data repository.

## Specific objectives:

- Enable RHSP researchers and statisticians correlate and analyze all study round data from one source for timely fulfillment of data requests.

- Develop a framework to perform this task repeatedly as additional survey, laboratory and clinical data is collected.

- Standardize data retrieval and ensure reproducible reports.

- Enforce data access controls to ensure data fidelity overtime.

- Proper documentation for easy variable definition and data auditing.

# Conceptual Framework

# Implementation(ETL)

**Source**
Organize FoxPro free tables into database containers for each round
➢Use FoxPro SEDNA upsizing wizard to load data into SQL staging tables within round specific databases
➢Freeze upsized FoxPro databases for consistence between FoxPro and SQL tables

**Staging**
➢Identify required domains (family, member, mobility, baseline, follow-up, partner, HIV, Viral load, CD4, Syphilis, BV )
➢Identify PII tables/columns for exclusion
➢Create data model for target tables

**Datamart**
➢Create data mart tables (exclude PII)
➢Create SSIS packages to load target data from the staging tables
  o Truncate staging tables once target data mart is loaded
➢Create views for routine or ad hoc queries and data visualization

# Staging tables

# SSIS packages

# Data mart tables and Views

# Business intelligence

# *Data Presentation*

- **Summarize** and visualize large amounts of historical data

- Only reveal **aggregate distributions** to mitigate risk of misuse

- Allow users to filter and represent data in **dynamic/interactive** ways

- Generate **new insights** that can lead to research hypotheses

- Leverage tableau public platform to enable access to broader research community
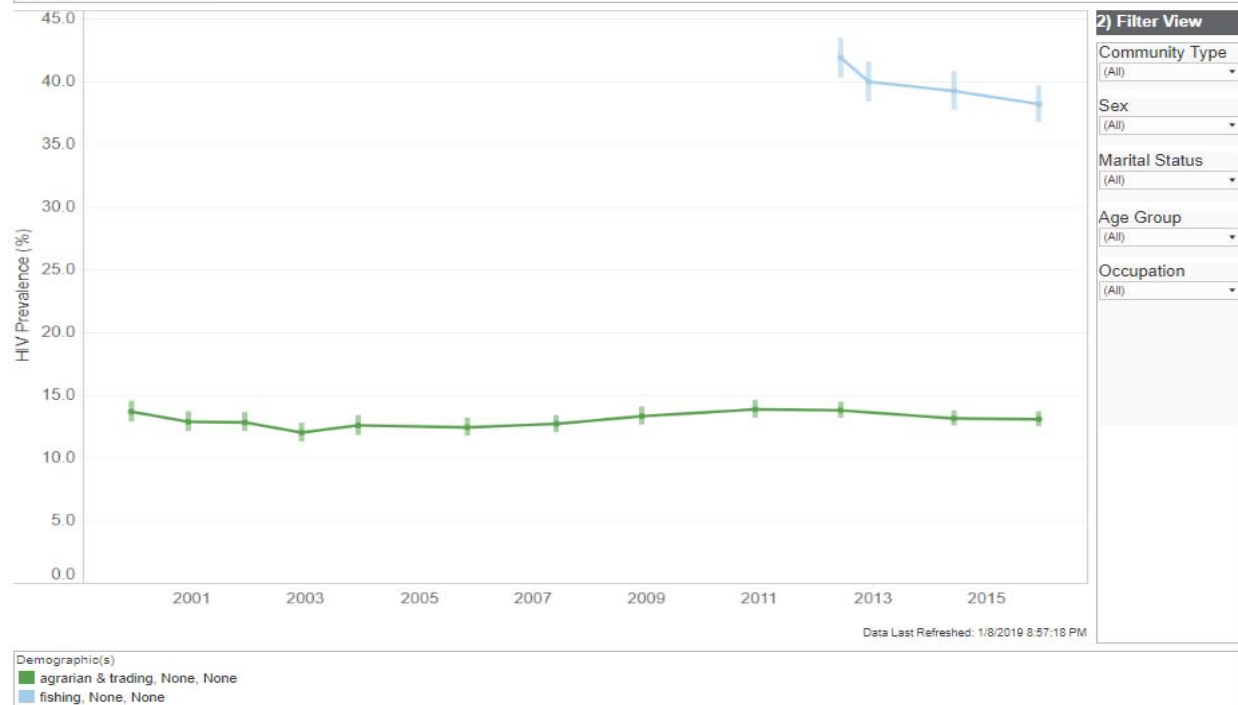


End Users (RHSP, JHU, NIAID)

# Current Dashboard

- Displays Prevalence trends from Rounds 6-17
- Separates Agrarian & Trading communities from Fishing communities
- Mirrors eligibility criteria from Nov 2017 NEJM article (https://www.nejm.org/doi/full/10.1056/NEJMoa1702150)
- Allows stratification by Sex, Age Group, Marital Status, and Occupation
- Currently adding Incidence, ART Coverage, and Male Circumcision Coverage
- Accessible for RHSP team via this link

# Results!

# RHSP Data Mart Solution

# RHSP Data Mart Solution

**Rakai Health Sciences Program**
*Improved Health Through Research*

**Initial Rollout: Feb 2015**

## Training Materials

- **10** Video Tutorials (~**150** minutes total)
- **2** User Manuals
- **2** User Test Plans
- **1** Training Plan

*Training resources to provide end users with the underlying knowledge and skillsets necessary to use the Data Mart*

## RHSP Data Dictionary

- **750** Pages
- **9** Columns of Information

*A compiled metadata repository filterable on database variables and their corresponding table locations/ questions/ answers*

**RHSP Data Mart**
*Source FoxPro tables (Rounds 1-15) were extracted, transformed, loaded and validated in the target SQL Server Data Mart*

## Data Model/ Design

- Data Mart Model
- Data Mart Design Doc.
- Physical Database Model
- Data Mart Table Links

*Data Mart physical and logical database models and their associated design documentations*

*Subset of queries derived from various scenarios and research use cases, including complex exclusion criteria*

## Library of Queries

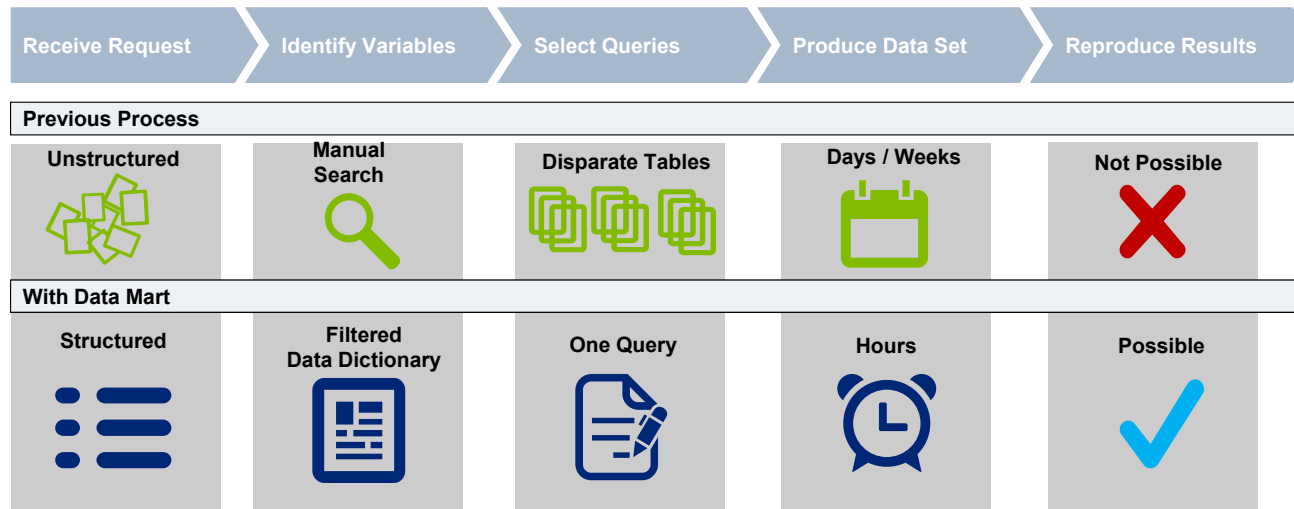- **34** Total sample Queries
- **91** ETL Scripts
- **5** Sample Use Cases

- **20** Tables
- **65,762** Subject Records
- **15+** Rounds of Research Data
- Track Changes Logs

# RHSP Data Mart Solution

| Receive Request | Identify Variables | Select Queries | Produce Data Set | Reproduce Results |
|---|---|---|---|---|

**Previous Process**

| Unstructured | Manual Search | Disparate Tables | Days / Weeks | Not Possible |
|---|---|---|---|---|

**With Data Mart**

| Structured | Filtered Data Dictionary | One Query | Hours | Possible |
|---|---|---|---|---|

Sample Query:

```
use RHSPDatamart
GO
select distinct Round,COUNT(study_id) baselines
from RHSPDatamart.dbo.Baseline1
where study_id<>' '
group by Round
order by Round
```

# *Conclusion*

- Moving towards more open information sharing with the wider community-allowing RHSP to be discovered through visual analytics and hopefully drive new collaborations and grants.

- Ready to harness new technologies to uncover meaningful insights in the data using predictive models, scenario analysis, and advanced analytics for timely and effective decision making.

- Legacy systems might have become outdated but setting up the infrastructure for modern systems requires high capital expenditure.

# Acknowledgments