# WHAT TOOLS AND METHODS DO WE NEED TO ENABLE SECONDARY USE OF ROUTINE DATA?

## Dr. Marcos Barreto

Associate Professor, Computer Science Department
Federal University of Bahia (UFBA), Salvador, Brazil

Farr Institute of Health Informatics Research, University College London
Newton International Fellow (The Royal Society, 2016-2018)

Ho Chi Minh City, Vietnam, 30 October to 2 November 2018
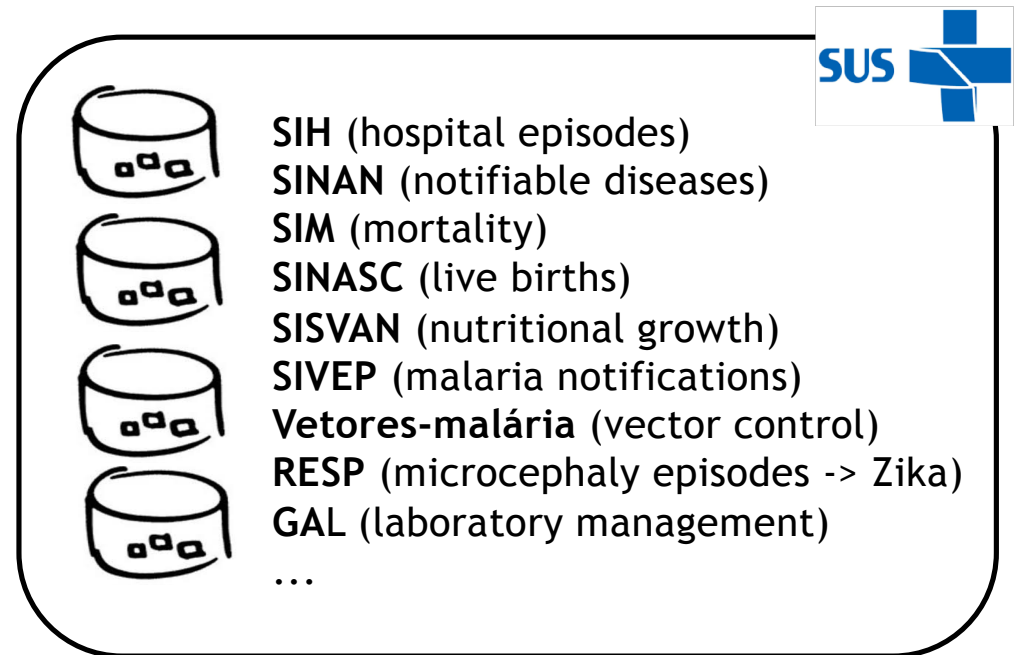
# Socioeconomic and health care routine data in Brazil

## Social programmes

✓ Targeted to poor and extremely poor families.
✓ Cadastro Único: central registrar for all programmes.

≅22 programmes

## Public health system (SUS)

✓ Big and complex public health system.
  - from primary care to specialised transplantations.

✓ Used by approximately 77% of the Brazilian population (164 million people).

**SIH** (hospital episodes)
**SINAN** (notifiable diseases)
**SIM** (mortality)
**SINASC** (live births)
**SISVAN** (nutritional growth)
**SIVEP** (malaria notifications)
**Vetores-malária** (vector control)
**RESP** (microcephaly episodes -> Zika)
**GAL** (laboratory management)
…

# Existing research platforms using these data

✓ **The 100 Million Cohort**



✓ **Zika surveillance (+ microcephaly)**



✓ **Malaria linkage & analytics**



✓ Baseline: CadastroÚnico, 2001-2015, 114 million individuals x 367 atributes.
✓ Cohort: baseline + Bolsa Família (cash transfers) + Housing (MCMV), 2001 – 2015, 400 million records, 3,000 attributes.
✓ Used by +20 projects assessing the effects of social programmes on health outcomes.

✓ Birth cohort, 2001 – 2030, ≅80 million records.
✓ Morbidity, mortality, socioeconomic and service data.
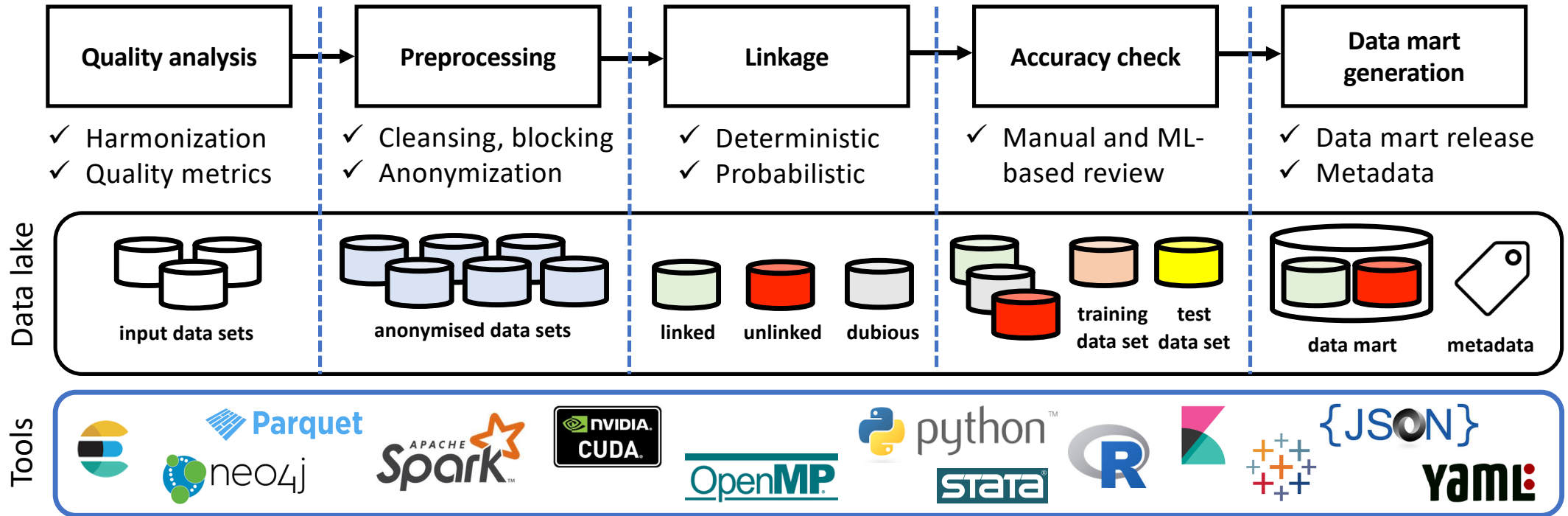✓ Focus on the triple epidemic (Zika, Dengue and Chikungunya) and health/educational outcomes related to microcephaly.

✓ Malaria episodes (>5 million records) + mortality + socioeconomic + climate data, 2000 – 2018.
✓ Focus on i) data aggregation and ii) epidemic forecasting.

# AtyImo – Data linkage platform

| Quality analysis | Preprocessing | Linkage | Accuracy check | Data mart generation |
|---|---|---|---|---|
| ✓ Harmonization<br>✓ Quality metrics | ✓ Cleansing, blocking<br>✓ Anonymization | ✓ Deterministic<br>✓ Probabilistic | ✓ Manual and ML-based review | ✓ Data mart release<br>✓ Metadata |

**Data lake**

input data sets — anonymised data sets — linked — unlinked — dubious — training data set — test data set — data mart — metadata

**Tools**

Parquet · neo4j · Spark · NVIDIA CUDA · OpenMP · python · stata · R · {JSON} · yaml

✓ Harmonization
✓ Data imputation

✓ Deep + Machine learning
✓ Statistical tools
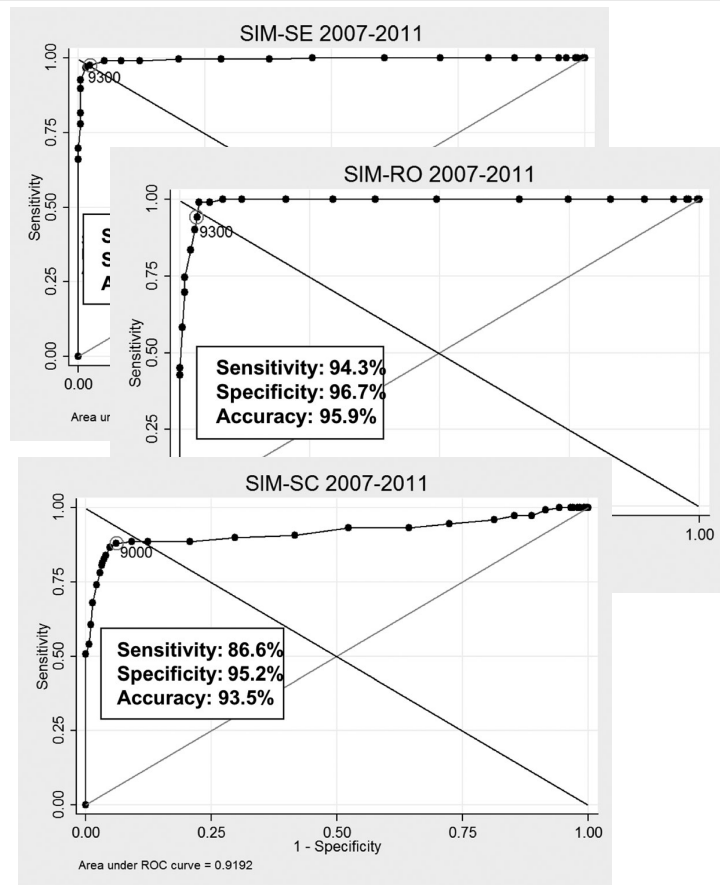
✓ Visual modelling, storyboards
✓ Geospatial data

**Data ingestion** → **Data analysis** → **Visual mining and report**

AtyImo
<www.atyimolab.ufba.br>

## Data analytics pipeline

# Example results



On the Accuracy and Scalability of Probabilistic Data Linkage Over the Brazilian 114 Million Cohort

IEEE JOURNAL OF BIOMEDICAL AND HEALTH INFORMATICS, VOL. 22, NO. 2, MARCH 2018

Robespierre Pita, Clícia Pinto, Samila Sena, Rosemeire Fiaccone, Leila Amorim, Sandra Reis, Mauricio L. Barreto, Spiros Denaxas, and Marcos Ennes Barreto

SIM-SE 2007-2011

SIM-RO 2007-2011

Sensitivity: 94.3%
Specificity: 96.7%
Accuracy: 95.9%

SIM-SC 2007-2011

Sensitivity: 86.6%
Specificity: 95.2%
Accuracy: 93.5%

Area under ROC curve = 0.9192

**Exploring hybrid parallel systems for probabilistic record linkage**

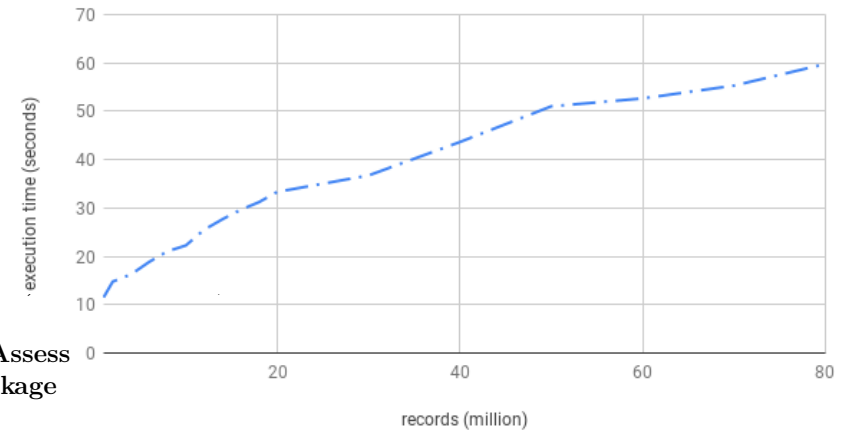Murilo Boratto[1] · Pedro Alonso[2] · Clicia Pinto[3] · Pedro Melo[3] · Marcos Barreto[3] · Spiros Denaxas[4]

**A Machine Learning Trainable Model to Assess the Accuracy of Probabilistic Record Linkage**

Robespierre Pita[1(✉)], Everton Mendonça[1], Sandra Reis[2], Marcos Barreto[1,3], and Spiros Denaxas[3]

Hybrid Execution Time

# Biggest challenges

✓ Technical:

    i)    <span style="color:red">Large volumes</span> of data hinder <span style="color:red">validation</span>.

    ii)   Data heterogeneity hinders the <span style="color:red">adoption of generalizable</span> machine learning <span style="color:red">methods</span>.

✓ Behavioural:

    ✓ How to <span style="color:red">conciliate</span> domain expertise (human) and powerful (but black-boxed) methods (machine) into <span style="color:red">accessible and effective</span> analytics platforms to support routine research and policy-making over linked health care data?

Thank You.

**Contact:**
**marcosb@ufba.br**
**www.atyimolab.ufba.br**