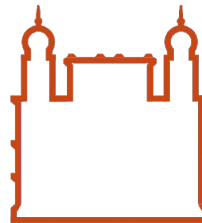


**cidacs**  
Centro de Integração de Dados  
e Conhecimentos para Saúde



Ministério da Saúde

**FIOCRUZ**

**Fundação Oswaldo Cruz**

Instituto Gonçalo Moniz

# Construction of Cohort 100M SINASC-SIM

# Summary

1. Overview of cohort 100M SINASC-SIM
2. Description of datasets
3. Pre-processing datasets
4. Harmonization
5. Record linkage

# Overview of cohort 100M SINASC-SIM



# Description of datasets

|                     | POP 100     |               | SIM        | SINASC     |
|---------------------|-------------|---------------|------------|------------|
|                     | CADU        | BOLSA FAMILIA |            |            |
| <b>Nº registers</b> | 114.008.179 | 27.376.582    | 17.829.111 | 44.485.274 |
| <b>Nº variables</b> | 259         | 21            | 138        | 86         |
| <b>Coverage</b>     | 2001-2015   | 2004-2015     | 2000-2015  | 2001-2015  |

# Pre-processing datasets

# Standardization: Integer variables

- Removal of special characters like:
  - (double) blank spaces;
  - \*, -, ', "
- The value becomes invalid when there is letters between the numbers

# Standardization: Integer variables

| Input     | Output |
|-----------|--------|
| 00123-000 | 123000 |
| 654*009   | 654009 |
| 01 0607   | 10607  |
| 01tgp09   | null   |

## Standardization: Date variables

- Removal of separators like:
  - -, /, space
- Conversion to date format
- The value becomes invalid when it cannot be converted to a known date format (yyyyMMdd, ddMMyyyy, etc...)



# Standardization: Date variables

| Input      | Output     |
|------------|------------|
| 01/01/2001 | 01-01-2001 |
| 11 12 2016 | 11-12-2006 |
| 2016-08-13 | 13-08-2016 |
| 40503799   | null       |

# Standardization: Categorical variables

- Checks if each category agrees with the dictionary.

# Standardization: Categorical variables (sex)

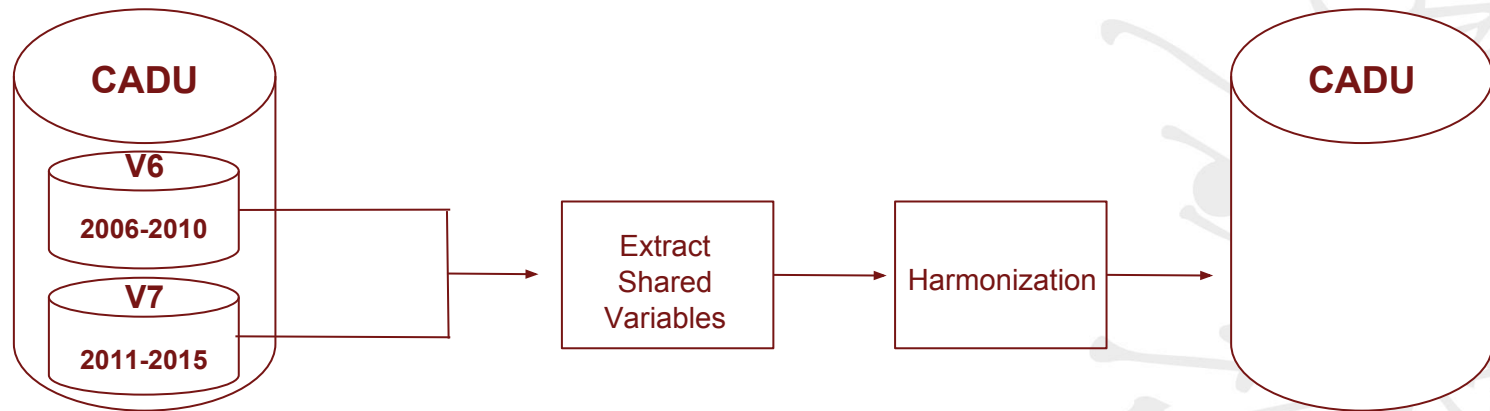
| Dictionary  |  |
|---|--|
| Gender:<br>M – Male<br>F – Female<br>I – Ignored          | Gender:<br>0 – Null<br>1 – Male<br>2 – Female<br>88 – Ignored<br>99 – Inconsistent |
| Input   | Output   |
| M   51000<br>F   54000<br>I   5800<br>O   2<br>null   700 | 0   700<br>1   51000<br>2   54000<br>88   5800<br>99   2                           |

# Pre-Processing: String type variables

- To improve the Linkage process
- To improve the analysis
- Depending on the dataset

| NAME             | NAME_CORR      |
|------------------|----------------|
| IGNORADO         | NULL           |
| JOAO SA*NTOS     | JOAO SANTOS    |
| MARIA2 DA SI8LVA | MARIA DA SILVA |

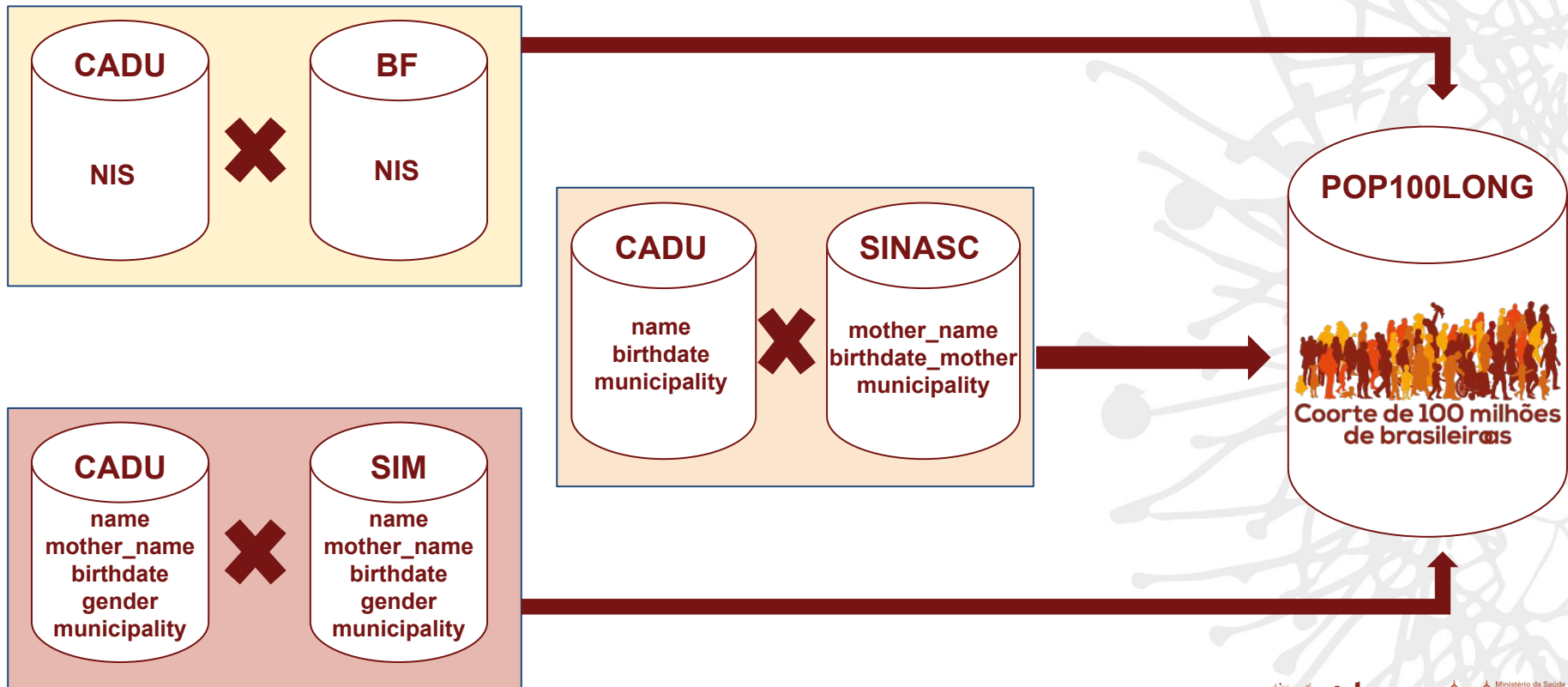
# Harmonization

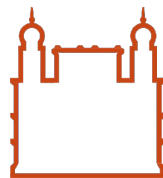


| GEN_V6 | GEN_V7 | GEN_H |
|--------|--------|-------|
| M      | 1      | 1     |
| F      | 2      | 2     |

# Record Linkage

# Record Linkage: POP100





Ministério da Saúde

FIOCRUZ

Fundação Oswaldo Cruz

Instituto Gonçalo Moniz

# Thanks !

#### COLABORADORES CIENTÍFICOS



#### APOIADORES



Parque Tecnológico da Bahia  
Rua Mundo, 121, Trogoy  
Salvador - Ba, CEP 41745-715



[www.cidacs.bahia.fiocruz.br](http://www.cidacs.bahia.fiocruz.br)  
[cidacs@bahia.fiocruz.br](mailto:cidacs@bahia.fiocruz.br)





# Applications

# Applications

From this dataset it is possible to know:

- Profile of people registered at CADU since birth to death;
- Factors (sociodemographic, maternal characteristics and childbirth) related to infant death;
- Living conditions of families whose children were born during the study period..

# Applications

To follow the impact of conditional cash transfer programme (Bolsa Familia) of some individual from the birth to the death.



# Standardization: Categorical variables

| Dictionary   |  |
|--|--|
| Method used. Values:<br>1 – Physical exam<br>2 – Another method<br>9 – Ignored | Method used. Values:<br>0 – Null<br>1 – Physical exam<br>2 – Another method<br>88 – Ignored<br>99 – Inconsistent |

| Input      | Output    |
|------------|-----------|
| 1   79000  | 0   100   |
| 2   1200   | 1   79000 |
| 3   97     | 2   1200  |
| 9   600    | 88   600  |
| null   100 | 99   97   |

# Standardization: Categorical variables which starts with 0

| Dictionary  |  |
|---|--|
| Schooling 2010. Valores:<br>0 – Without Schooling;<br>1 – Fundamental I (1ª a 4ª year);<br>2 – Fundamental II (5ª a 8ª year);<br>3 – Secondary (Old second degree);<br>4 – Incomplete undergraduated;<br>5 – Ungraduated;<br>9 – Ignored. | Schooling 2010. Valores:<br>0 – Null<br>1 – Without Schooling<br>2 – Fundamental I (1ª a 4ª série)<br>3 – Fundamental II (5ª a 8ª série)<br>4 – Secondary (Old second degree)<br>5 – Incomplete undergraduated;<br>6 – Ungraduated<br>88 – Ignored<br>99 – Inconsistent. |

# Standardization: Integer variables

| Dictionary |                   |
|------------|-------------------|
| Birth hour | Birth hour (HHMM) |

| Input | Output |
|-------|--------|
| 08:59 | 0859   |
| 09/36 | 0936   |
| 085:1 | null   |

# Standardization: Binary Categorical variables

| Dicionário                                      |  |
|---|--|
| Status of new DO variable:<br>1 – Yes<br>0 – No | Status of new DO variable:<br>0 – Null<br>1 – Yes<br>2 – No<br>99 – Inconsistent |

| Entrada                                       | Saída                                       |
|---|---|
| 1   45000<br>0   6700<br>3   22<br>null   700 | 0   700<br>1   45000<br>2   6700<br>99   22 |

# Correction on string type variables

| Resultados                    |         |
|-------------------------------|---------|
| Recovered names of children   | 0       |
| Unrecovered names of children | 419.557 |
| Recovered names of mothers    | 5.228   |
| Unrecovered names of mothers  | 1.932   |