# Cell Phone Data for Social Impact

## MIT CDR Data Training

**JAN 2020**

# Agenda

| | |
|---|---|
| 11:00 - 11:30 AM | Preparation |
| 11:30 AM - 12:00 PM | Introduction to Spark |
| 12:00 - 12:15 PM | Break |
| 12:15 - 12:30 PM | Data Preparation |
| 12:30 - 1:15 PM | Data Processing and Exploration |
| 1:15 - 1:30 PM | Closing & Conclusion |

# Goals & Objectives

1. Understand some basic tools for large-scale, big data processing

2. Wrangle and ready CDR data for processing

3. Learn the basics of Apache Spark

4. Explore the use of Apache Zeppelin for integrated analytics

5. Work collaboratively in a secure analytics environment

# Tools & Frameworks

| |
|---|
| Parallel/Distributed Processing |
| Storage and Computing |
| Package Management |
| Analytics and Data Processing |
| Programmatic API |

# X
# Setting Up

## X  Setting Up

1. Install <u>FoxyProxy</u>
2. Download the <u>FoxyProxy settings file</u>
3. Import the settings file to Foxy Proxy
4. Select the proxy setting:
   1. `Use proxy emr-socks-proxy for all URLs`
5. Open your Terminal/Shell and enter the following:
   1. `ssh -i ~/.ssh/<private-key> -ND 8157 <username>@34.223.103.224`
6. Return to your browser
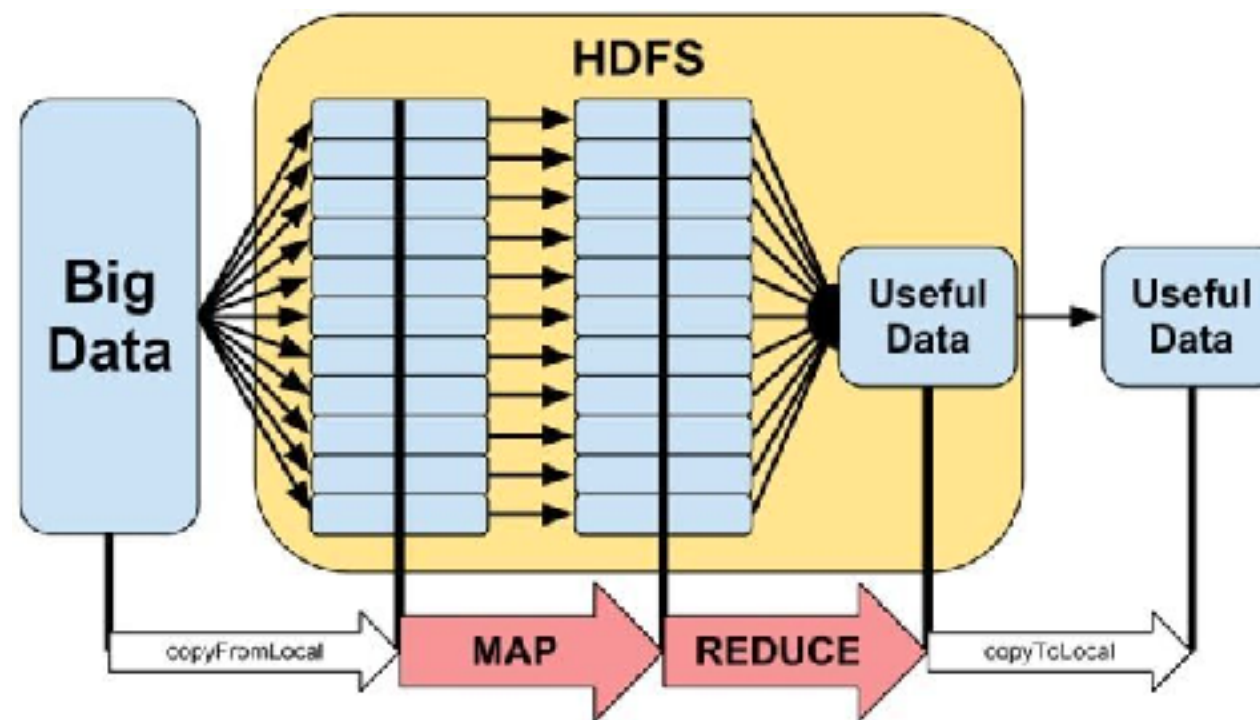
## X  Setting Up: Why?

Our instance is not accessible via the public internet. This is good. In order to access our analytics environment, you need secure (authenticated and authorized).

We have two basic layers of security: SSH and AWS IAMs

We have control over who has access to CDRs and know what they are doing with it!
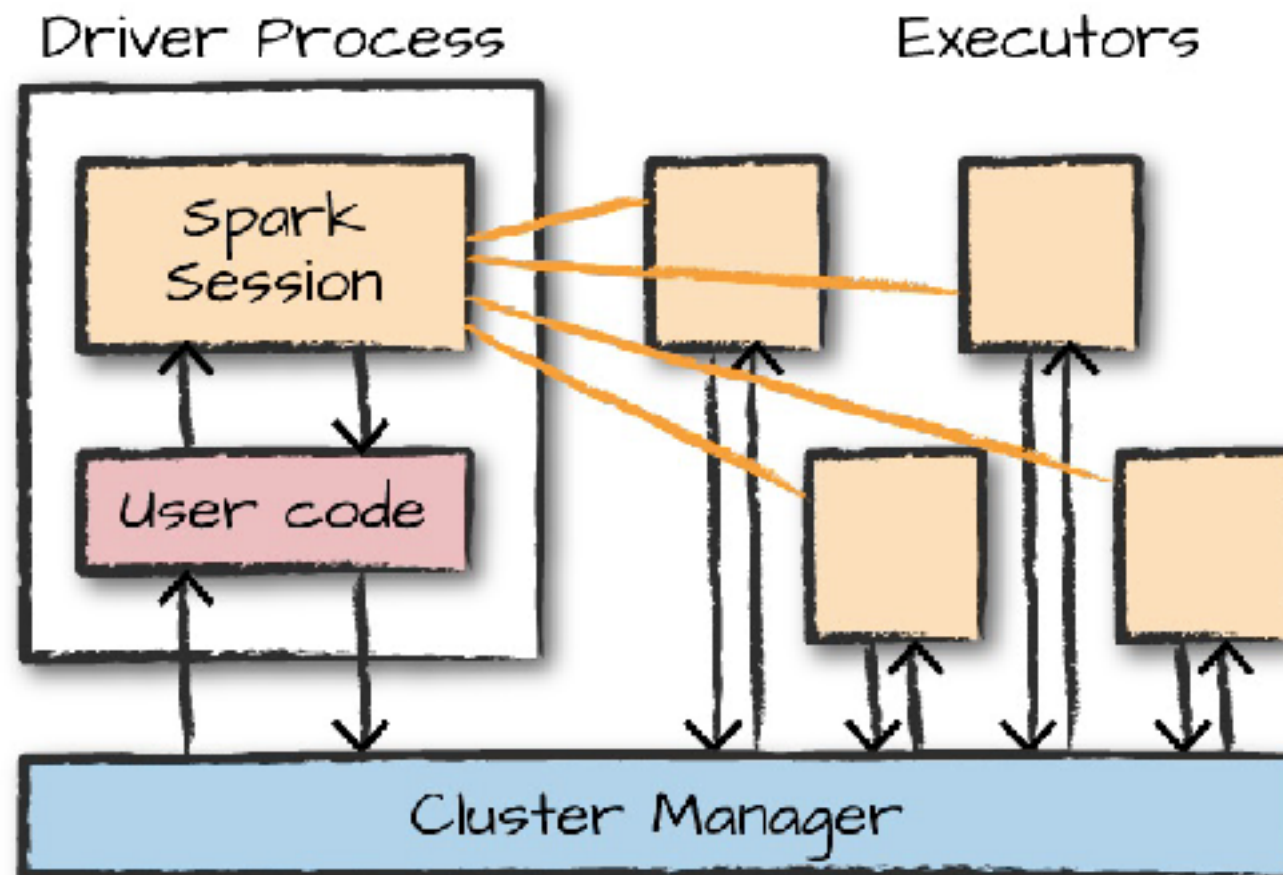
MITGOV/LAB

CIVIC DESIGN
DATA LAB

DIRECTORATE OF SCIENCE
TECHNOLOGY & INNOVATION

africell

# I
# Introduction to Apache Spark

**I** **What is Spark?**



**a better Hadoop?**

# What is Spark?



Source: Chamber, Bill et al. 2017. Spark: The Definitive Guide, pg 21

- Spark processes in-memory and can run onto various filesystems (HDFS, S3, RDBMs, etc).

- Very powerful APIs and abstraction.

- Spark is lazy.

**I** **What is Spark?**

Follow on Zeppelin

## I    What is Spark: Dataframe

Spreadsheet on
a single machine

Table or Data Frame
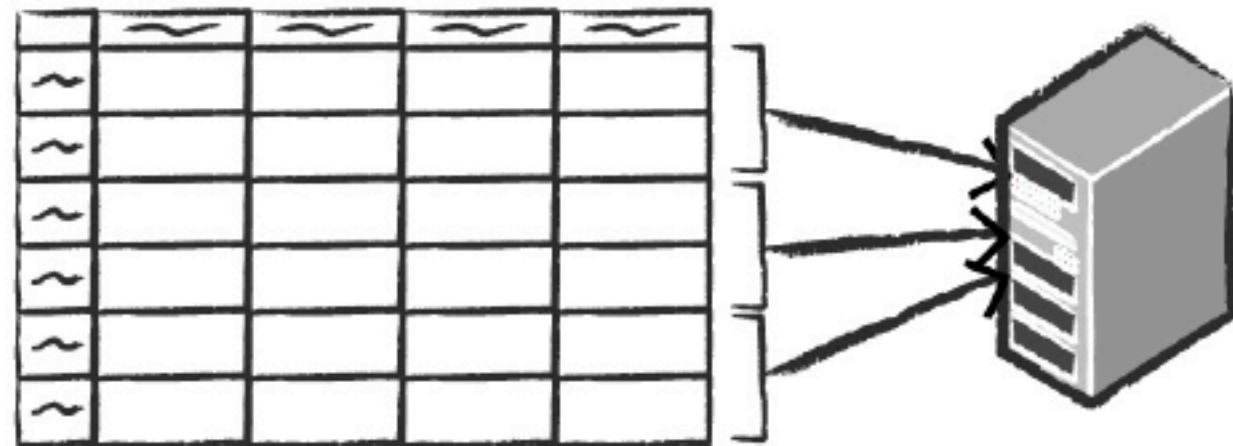partitioned across servers
in a data center

Figure 2-3. Distributed versus single-machine analysis

# I   What is Spark: Transformations
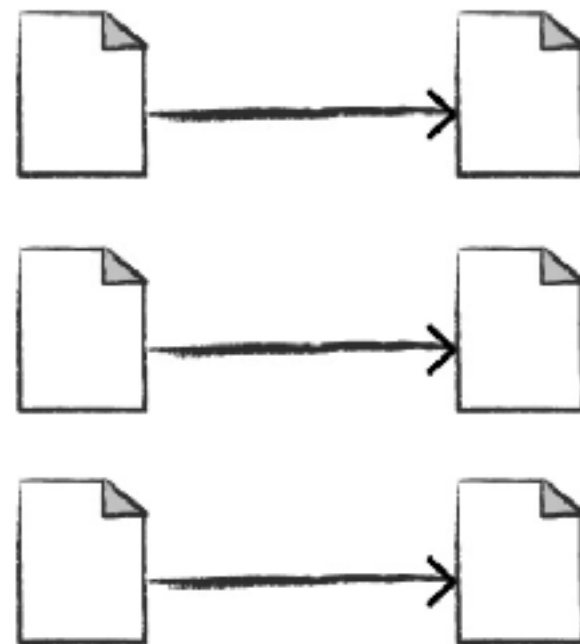


Narrow transformations
I to I

Figure 2-4. A narrow dependency
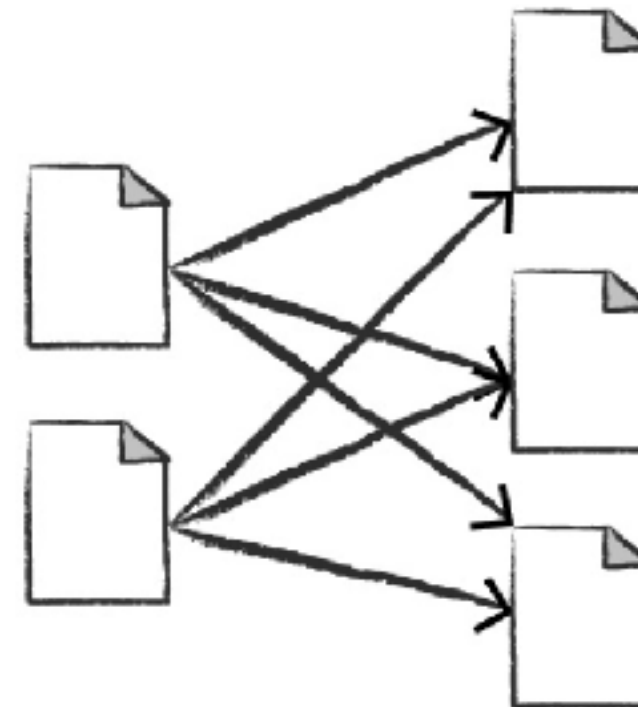
Wide transformations
(shuffles) I to N

Figure 2-5. A wide dependency

*Source: Chamber, Bill et al. 2017. Spark: The Definitive Guide, pg 27-28*

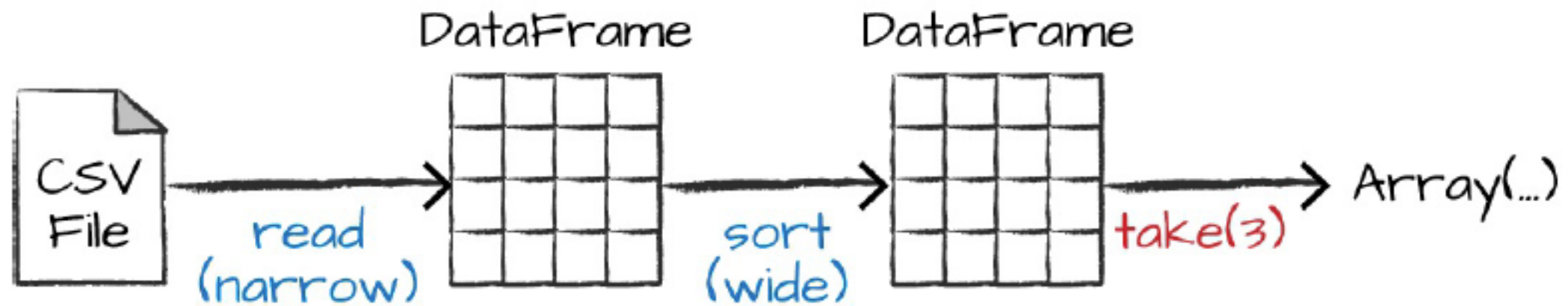# What is Spark: Transformations



Figure 2-8. Reading, sorting, and collecting a DataFrame

*Source: Chamber, Bill et al. 2017. Spark: The Definitive Guide, pg 31*

# What is Spark: Transformations



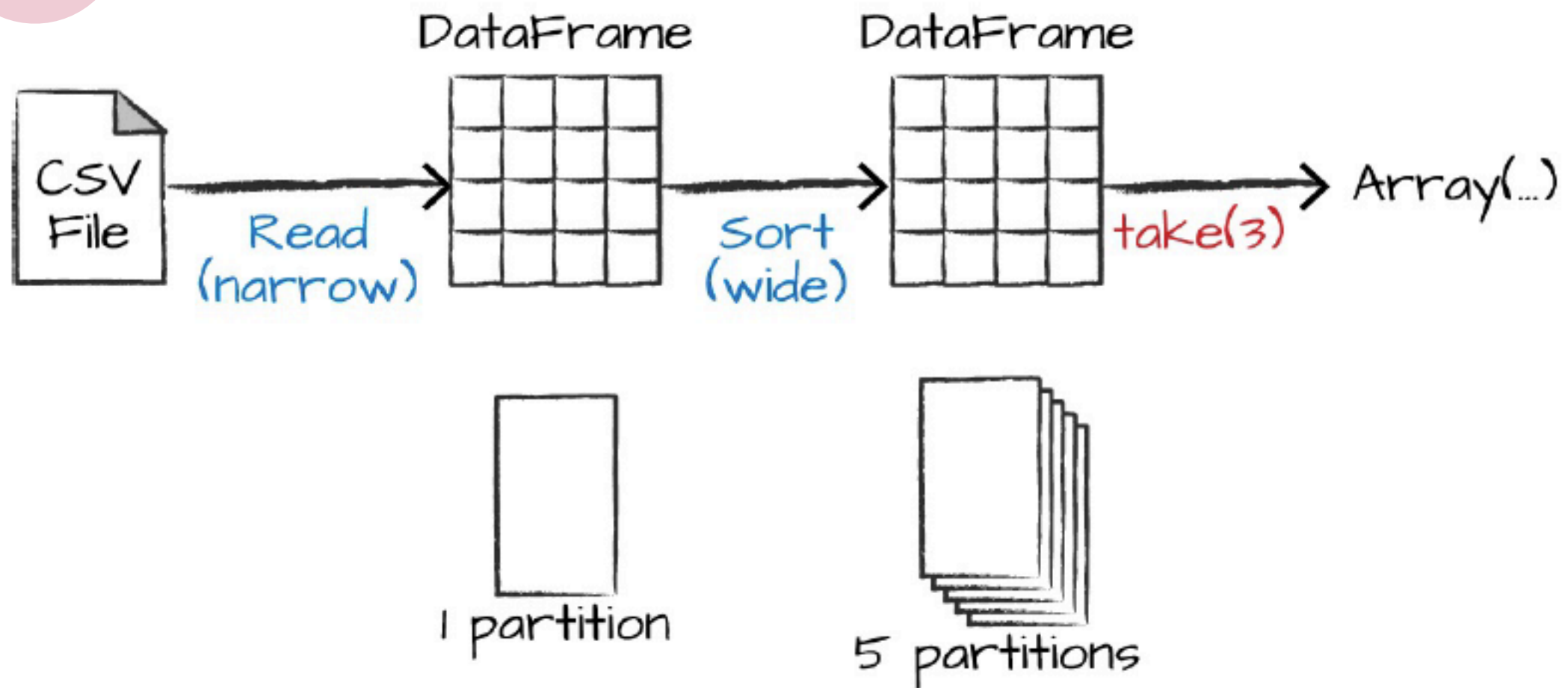Figure 2-9. The process of logical and physical DataFrame manipulation

Source: Chamber, Bill et al. 2017. Spark: The Definitive Guide, pg 32

# What is Spark: Transformations



Figure 5-2. Different kinds of transformations

- Remove columns or rows
- Transform a row into a column or a column into a row
- Add rows or columns
- Sort data by values in rows

*Source: Chamber, Bill et al. 2017. Spark: The Definitive Guide, pg 32*
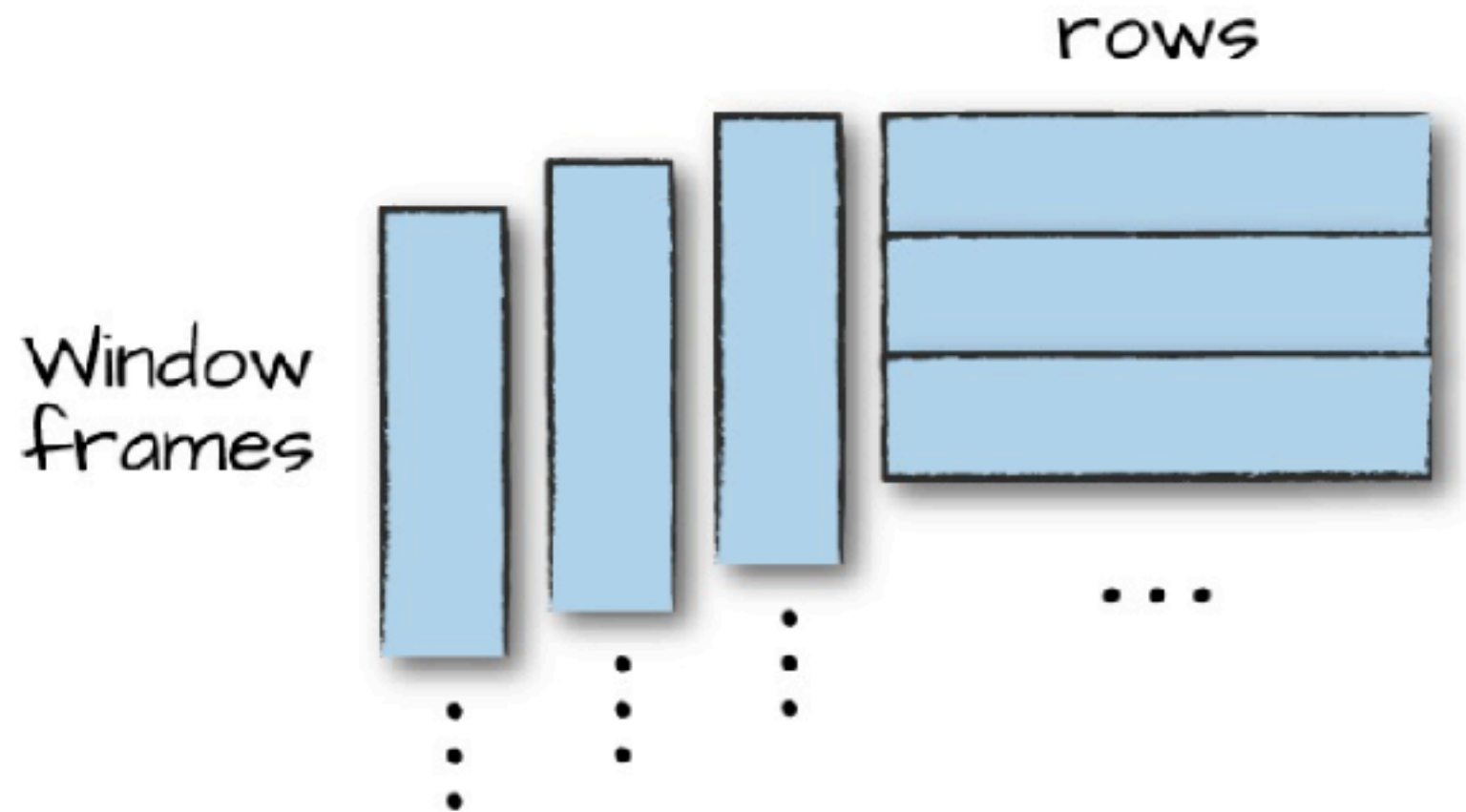
# **I** **What is Spark: Windows**



Figure 7-1. Visualizing window functions

Source: Chamber, Bill et al. 2017. Spark: The Definitive Guide, pg 135

# II
# Prepping Data For Processing

# Prepping Data for Processing



**Figure 1:** *Basic Big Data Analytics Processing*

**II**  **Prepping Data for Processing**

Our analyses require other types of data:

1. **shapefiles** for districts, chiefdoms, sections

2. antenna and site **meta-data**

We need to be intentional in our processing of this data. The decision we make here will impact the rest of our analysis. Time should be spent cleaning and understanding these datasets.

# Follow on Zeppelin

# III

# Creating Aggregates

We will review to types of aggregation task in Spark:

1. **custom aggregations**
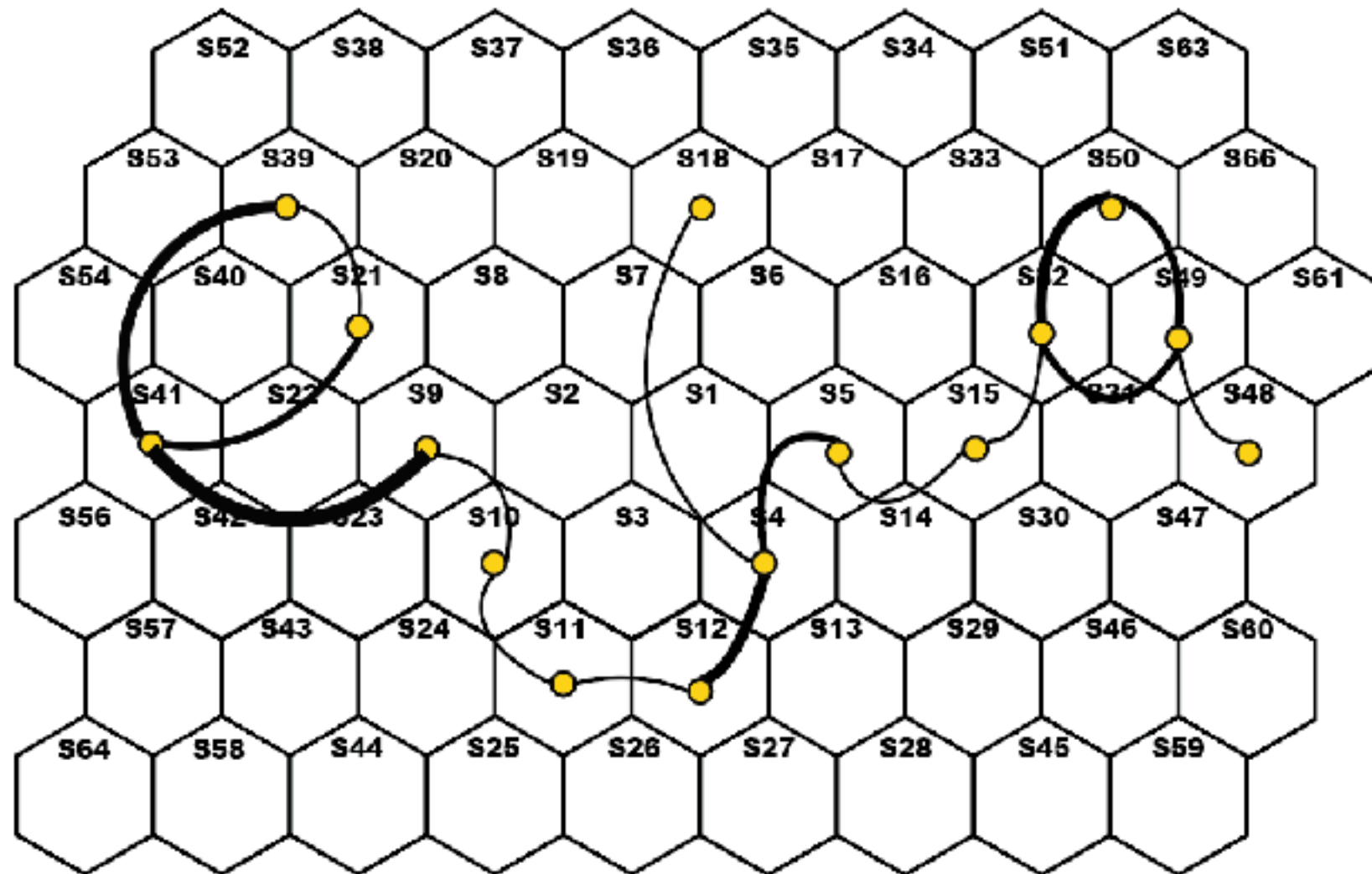
    1. generating stops and journeys

2. **off-the-shelf aggregations**

    1. running <u>Flowminder</u> SQL Queries

# III Creating Aggregates: stops & journeys

**III** **Creating Aggregates: stops & journeys**

**Follow on Zeppelin**

# IV
# Closing

**Closing: Privacy**

| 01 | **Identity Disclosure** | • Occurs when an individual is linked to a particular record in a released table |
|----|------------------------|------------------------------------------------------------------------------|
| 02 | **Attribute Disclosure** | • Occurs when new information about some individual is revealed and the released data makes it possible to infer the characteristics of an individual more accurately than it would be possible before release. |

For a comprehensive overview, see *Protecting User Data and Privacy*

**Privacy Preserving Methods**

- De-identification
- Adding and Hiding Events
- Anonymization
- Sampling
- Data Suppression/Swapping
- Perturbation

1 Location obfuscation

2 Statistical perturbation

For a comprehensive overview, see *Protecting User Data and Privacy*

MIT GOV/LAB    CIVIC DATA DESIGN LAB    DIRECTORATE OF SCIENCE TECHNOLOGY & INNOVATION    africell

# IV  Closing: Privacy

**Can the Data Analysis partner ensure:**

| 1 Effective Anonymization | 2 Appropriate Storage and Access | 3 High Standard of Ethical Use |
|---|---|---|
| - All Personally Identifiable Information and sensitive information are removed from dataset | - Anonymized CDR data is kept in a secure location accessible only by those with proper authorization. | - No attempt will be made to use external data and a priori information of identify individuals |
| - All analysis should protect the privacy of all individuals by utilizing methods that preserve a given level of anonymity | - Anonymized CDR should not be moved outside of the chosen storage environment. | - Requires data analysis partner to report any exposure of sensitive information |
| | | - Research outputs should not be capable of re-identification if combined with other data. |

For a comprehensive overview, see *Protecting User Data and Privacy*

MITGOV/LAB  CIVIC DATA DESIGN LAB  DIRECTORATE OF SCIENCE TECHNOLOGY & INNOVATION  africell

**IV**   **Closing: Much More**

- CDR can be use to model social networks. Spatial and temporal interactions among subscribers

- Can be overplayed with other data: census, business, survey.

- Apache Spark can handle this! Support for both streaming and batch analytics. Integrations pipelines for large-scale Machine Learning.

# Thank You!

**AX** References