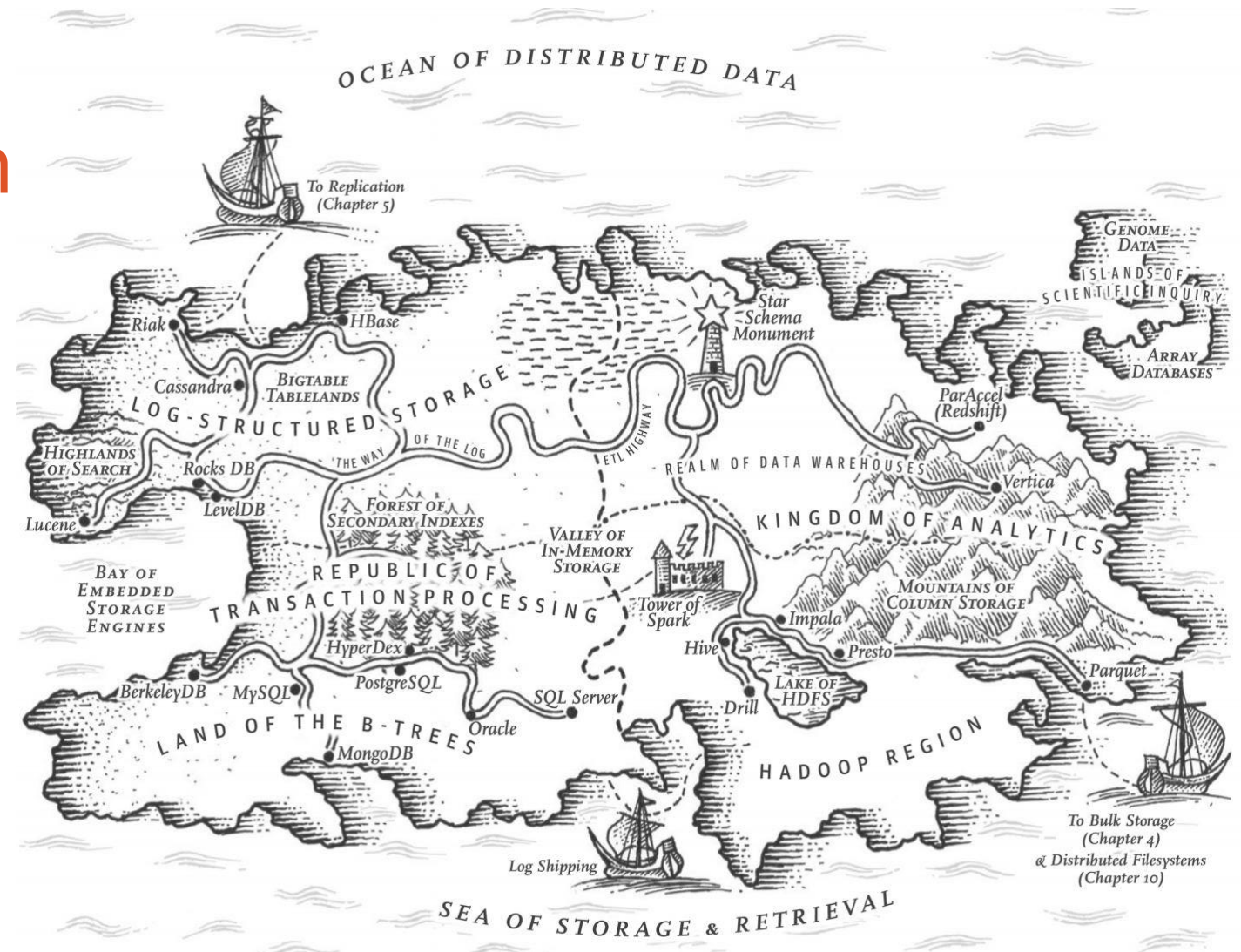A gentle introduction to

# Apache Spark, Databricks and Delta Lake

# Personal Introduction

- Senior Consultant at Avanade
- MSc in Financial & Management Engineering
- Microsoft Certified Professional
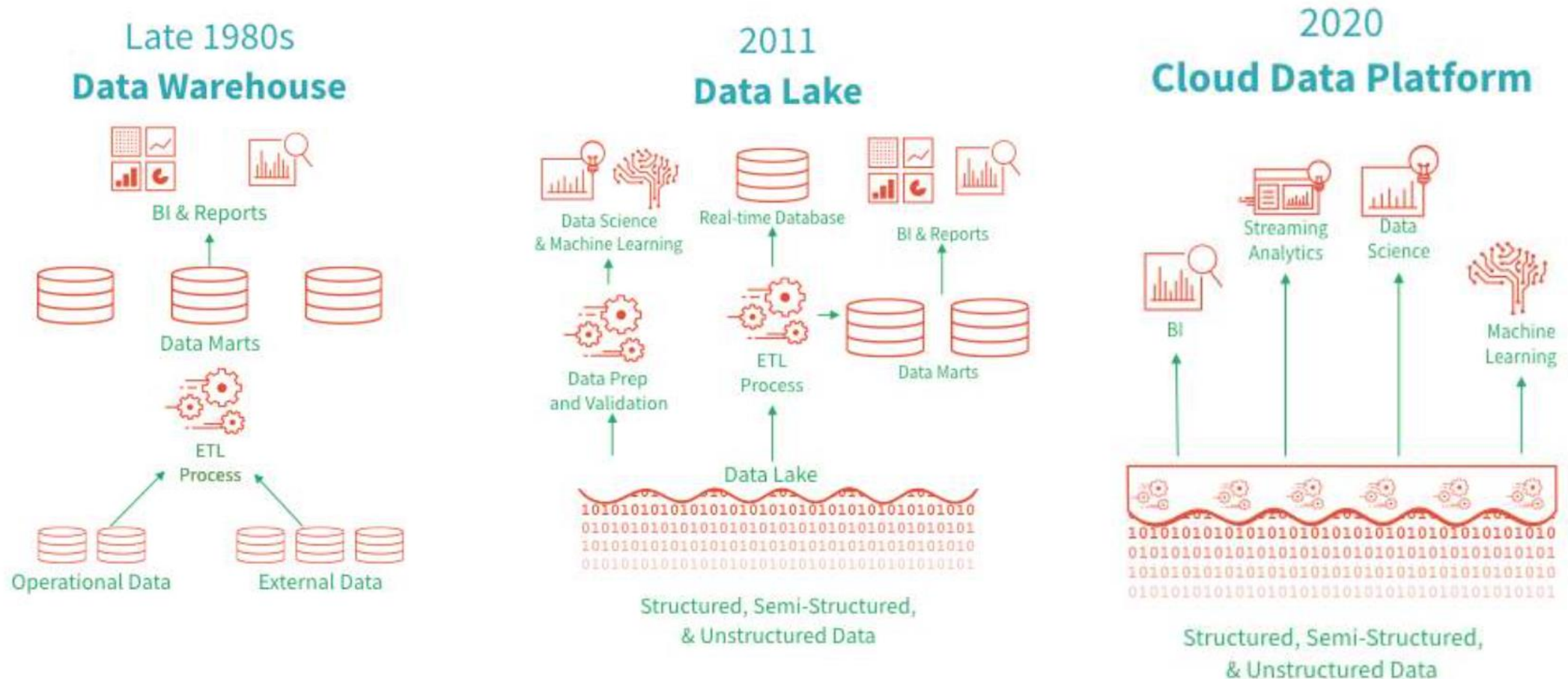- 12 years long journey, in the vast continent of Data



*Designing Data-Intensive Applications by Martin Kleppmann*

# Agenda

- The big picture

- A peak into Apache Spark using Databricks

- Delta Lake

- Where do I go from here?

# The evolution of Data Management



**Late 1980s**
**Data Warehouse**

BI & Reports

Data Marts

ETL Process

Operational Data    External Data

**2011**
**Data Lake**

Data Science & Machine Learning    Real-time Database    BI & Reports

Data Prep and Validation    ETL Process    Data Marts

Data Lake

1010101010101010101010101010101010101010
0101010101010101010101010101010101010101
1010101010101010101010101010101010101010
0101010101010101010101010101010101010101

Structured, Semi-Structured, & Unstructured Data

**2020**
**Cloud Data Platform**

Streaming Analytics    Data Science

BI    Machine Learning

1010101010101010101010101010101010101010
0101010101010101010101010101010101010101
1010101010101010101010101010101010101010
0101010101010101010101010101010101010101

Structured, Semi-Structured, & Unstructured Data

# A brief history of Apache Spark

- Began in 2009 at UC Berkeley as a research project

- At that time Hadoop was the dominant paradigm

- Databricks was founded in 2013

- Spark 1.0 released in 2014

- Spark 2.0 released in 2016

- Spark 3.0 released in 2020
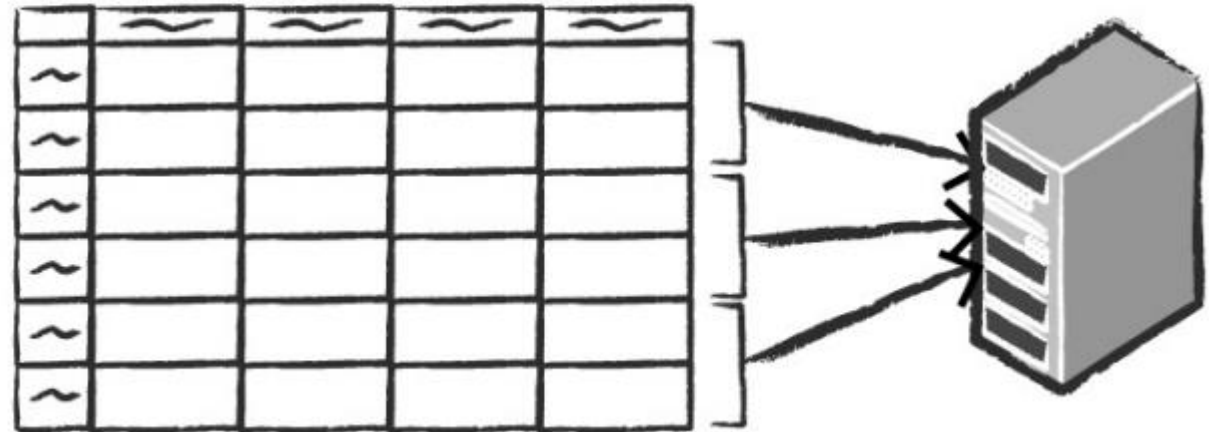
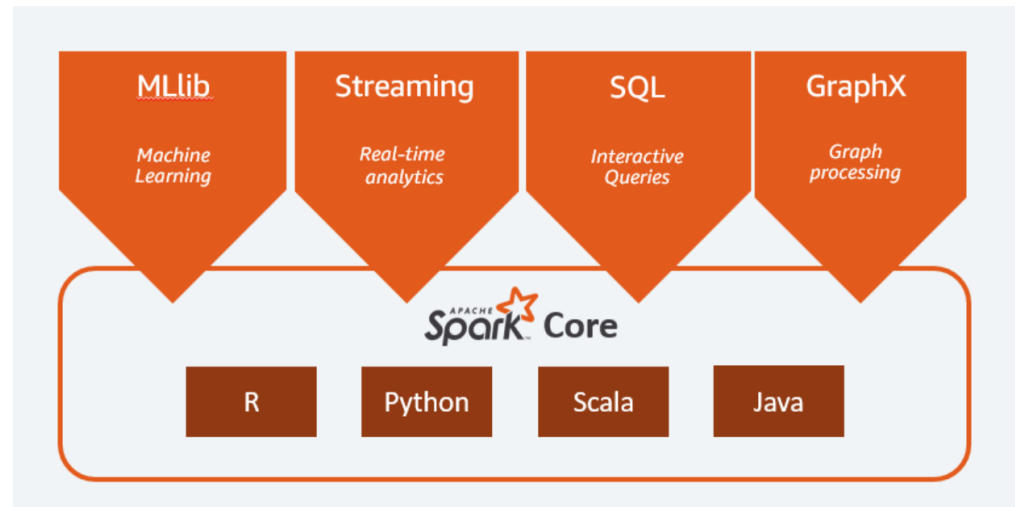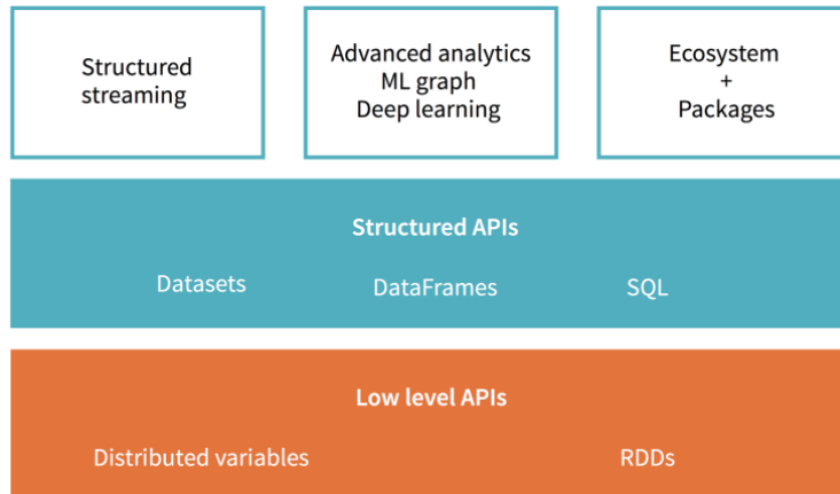What is Spark and how does it work?
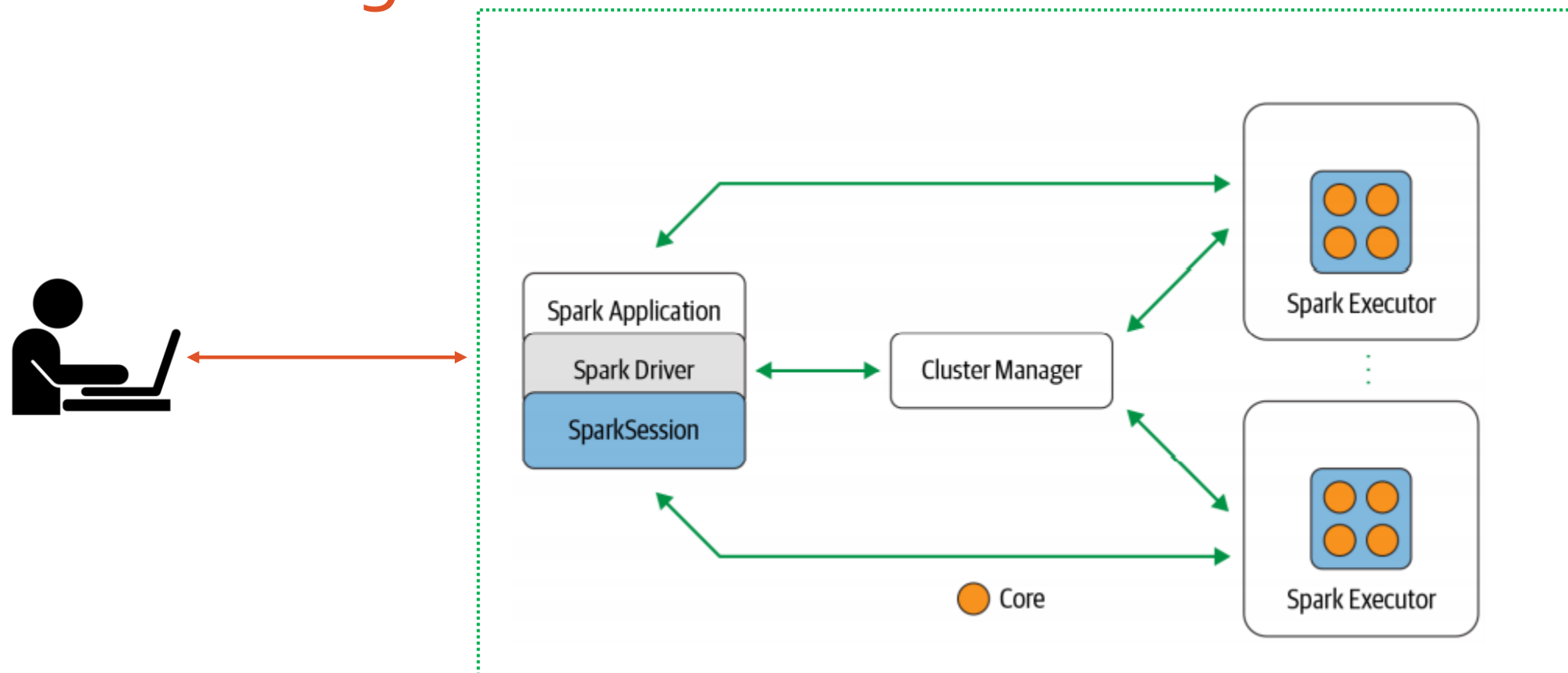
# A simple example



Spreadsheet on
a single machine

Table or Data Frame
partitioned across servers
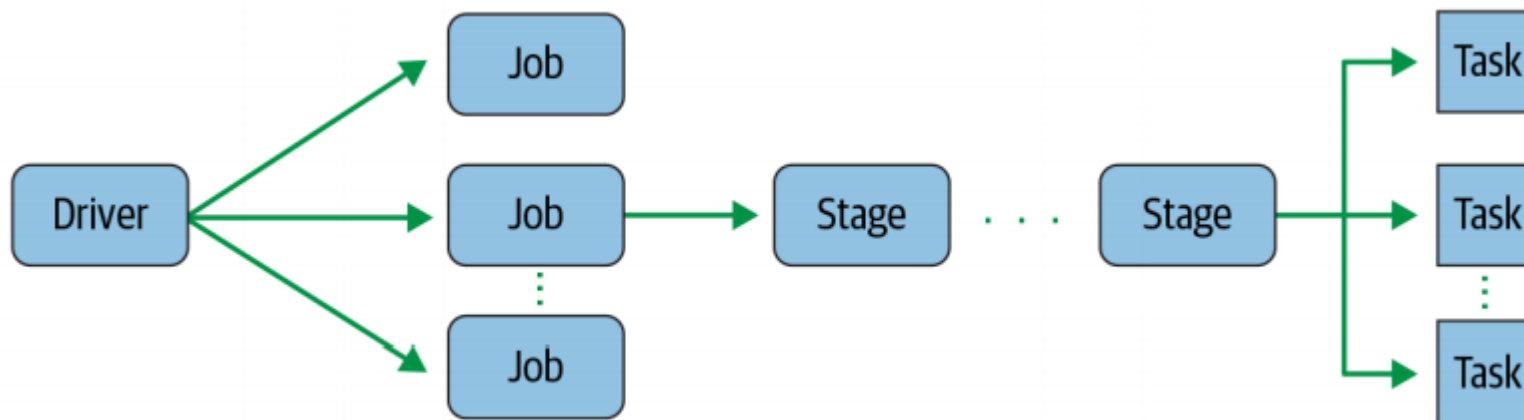in a data center

# How do we talk to Spark?

# Distributing work

# Distributing work

# Spark stages, jobs and tasks

# Transformations, Actions and Lazy Evaluation



```python
# In Python
>>> strings = spark.read.text("../README.md")
>>> filtered = strings.filter(strings.value.contains("Spark"))
>>> filtered.count()
20

// In Scala
scala> import org.apache.spark.sql.functions._
scala> val strings = spark.read.text("../README.md")
scala> val filtered = strings.filter(col("value").contains("Spark"))
scala> filtered.count()
res5: Long = 20
```
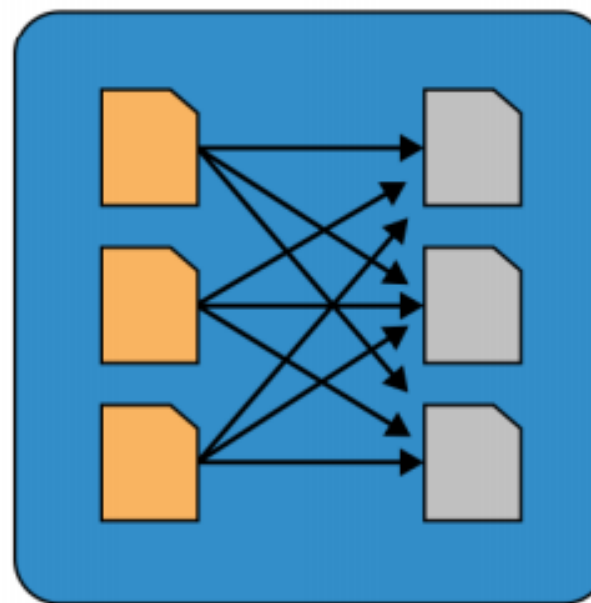
# Transformations, Actions and Lazy Evaluation

# An example on narrow transformations:

*"Select data where age = 37"*

**Node 1**

| Name | Age | City |
|------|-----|------|
| Arnold | 37 | Amsterdam |
| Mohamed | 37 | London |
| John | 25 | Athens |

**Node 1**

| Name | Age | City |
|------|-----|------|
| Arnold | 37 | Amsterdam |
| Mohamed | 37 | London |

**Node 2**

| Name | Age | City |
|------|-----|------|
| Lara | 37 | New York |
| George | 31 | London |
| Seif | 45 | Cairo |

**Node 2**

| Name | Age | City |
|------|-----|------|
| Lara | 37 | New York |

**Node 3**

| Name | Age | City |
|------|-----|------|
| Ankur | 37 | Mumbai |
| Jack | 67 | London |
| Lian | 24 | Beijing |

**Node 3**

| Name | Age | City |
|------|-----|------|
| Ankur | 37 | Mumbai |

**Final Result**

| Name | Age | City |
|------|-----|------|
| Arnold | 37 | Amsterdam |
| Mohamed | 37 | London |
| Lara | 37 | New York |
| Ankur | 37 | Mumbai |

# Delta Lake

# What is Delta Lake?

- Specifically designed to work with Apache Spark

- ACID compliant

- Streaming and batch unification

- Schema enforcement

- Time travel

- Upserts and deletes

# Example Delta Lake architecture

# How can I use Spark?

- The 'unmanaged way'

- The 'managed way'

# Demo

# Source: Databricks Documentation

# Where do I go from here

- https://www.oreilly.com/library/view/learning-spark-2nd/9781492050032/



- https://academy.databricks.com/pathway/how-to-build-a-cloud-data-platform



- https://docs.databricks.com/index.html

# Thank you!

Any questions?